

## Дополнительные материалы

Материалы > Анализ и обработка данных

### Тест Литтла для выявления MCAR

Для того чтобы количественно определить полностью случайным ли образом сформированы пропущенные значения (MCAR), существует тест Литтла (Little's test). Этот критерий был предложен в 1988 году Родериком Литтлом в качестве единого критерия оценки случайности пропущенных значений в многомерных *количественных* (!) данных.

Нулевая гипотеза этого критерия утверждает, что пропуски полностью случайны (MCAR). Альтернативная гипотеза говорит о том, что пропуски зависят от наблюдаемых значений (MAR).

Создадим новый блокнот на R в Google Colab и воспользуемся соответствующей функцией одной из библиотек этого языка.

Откроем блокнот с кодом на R

### Датасет airquality

Вначале попрактикуемся на встроенном в R датасете airquality.

```
1 # импортируем библиотеку datasets
2 library(datasets)
```

Посмотрим на первые строки датафрейма с помощью **функции head()**.

```
1 head(airquality)
```

A data.frame: 6 × 6

	Ozone	Solar.R	Wind	Temp	Month	Day
	<int>	<int>	<dbl>	<int>	<int>	<int>
1	41	190	7.4	67	5	1
2	36	118	8.0	72	5	2
3	12	149	12.6	74	5	3
4	18	313	11.5	62	5	4
5	NA	NA	14.3	56	5	5

5	NA	NA	14.9	66	5	5
6	28	NA	14.9	66	5	6

Оценим размерность.

```
1 # выведем общее количество строк и столбцов
2 dim(airquality)
```

```
1 153 6
```

Посмотрим на общие статистические показатели с помощью функции **summary()**.

```
1 summary(airquality)
```

```
      Ozone      Solar.R      Wind      Temp
Min.   : 1.00   Min.   : 7.0   Min.   : 1.700   Min.   :56.00
1st Qu.: 18.00   1st Qu.:115.8   1st Qu.: 7.400   1st Qu.:72.00
Median : 31.50   Median :205.0   Median : 9.700   Median :79.00
Mean   : 42.13   Mean   :185.9   Mean   : 9.958   Mean   :77.88
3rd Qu.: 63.25   3rd Qu.:258.8   3rd Qu.:11.500   3rd Qu.:85.00
Max.   :168.00   Max.   :334.0   Max.   :20.700   Max.   :97.00
NA's   :37      NA's   :7

      Month      Day
Min.   :5.000   Min.   : 1.0
1st Qu.:6.000   1st Qu.: 8.0
Median :7.000   Median :16.0
Mean   :6.993   Mean   :15.8
3rd Qu.:8.000   3rd Qu.:23.0
Max.   :9.000   Max.   :31.0
```

Теперь установим и импортируем библиотеку **nanian**<sup>4</sup>, которая и содержит необходимый нам тест Литтла.

```
1 # установим библиотеку
2 install.packages('nanian')
```

```
1 # импортируем ее
2 library(nanian)
```

Посмотрим на абсолютное количество и процент пропусков в каждом из столбцов.

```
1 miss_var_summary(airquality)
```

A tibble: 6 × 3

variable	n_miss	pct_miss
<chr>	<int>	<dbl>
Ozone	37	24.183007
Solar.R	7	4.575163
Wind	0	0.000000
Temp	0	0.000000
Month	0	0.000000
Day	0	0.000000

Перейдем к статистическому тесту.

```
1 # для теста Литтла используем функцию mcar_test()
2 mcar_test(airquality)
```

A tibble: 1 × 4

statistic	df	p.value	missing.patterns
<dbl>	<dbl>	<dbl>	<int>
35.10613	14	0.001417781	4

Вероятность (p-value) очень мала, ниже порогового значения, например, в пять или даже один процент, и мы можем отвергнуть нулевую гипотезу о полностью случайных пропусках.

Кроме этого, **функция mcar\_test()** сообщает о четырех выявленных в пропущенных данных паттернах.

Вернемся к датасету, с которым работали до сих пор.

## Датасет «Титаник»

Подгрузим файл train.csv в сессионное хранилище блокнота на R.

train.csv

[Скачать](#)

Импортируем его с помощью **функции read.csv()**.

```
1 titanic = read.csv('/content/train.csv', na.strings = c("NA", "NaN", ""))
2 head(titanic)
```

A data.frame: 6 × 12

PassengerId	Survived	Pclass	Name	Sex	Age	SibSp	Parch	Ticket	Fare	Cabin	Embarked	
<int>	<int>	<int>	<chr>	<chr>	<dbl>	<int>	<int>	<chr>	<dbl>	<chr>	<chr>	
1	1	0	3	Braund, Mr. Owen Harris	male	22	1	0	A/5 21171	7.2500	NA	S
2	2	1	1	Cumings, Mrs. John Bradley (Florence Briggs Thayer)	female	38	1	0	PC 17599	71.2833	C85	C
3	3	1	3	Heikkinen, Miss. Laina	female	26	0	0	STON/O2. 3101282	7.9250	NA	S
4	4	1	1	Futrelle, Mrs. Jacques Heath (Lily May Peel)	female	35	1	0	113803	53.1000	C123	S
5	5	0	3	Allen, Mr. William Henry	male	35	0	0	373450	8.0500	NA	S
6	6	0	3	Moran, Mr. James	male	NA	0	0	330877	8.4583	NA	Q

Обратите внимание, в параметр na.strings мы передали вектор с возможными вариантами записи пропущенных значений. Без этого параметра часть пропусков не была бы учтена.

Посмотрим на количество и процент пропусков в каждом столбце.

```
1 miss_var_summary(titanic)
```

A tibble: 12 × 3

variable	n_miss	pct_miss
<chr>	<int>	<dbl>
Cabin	687	77.1043771
Age	177	19.8653199
Embarked	2	0.2244669
PassengerId	0	0.0000000
Survived	0	0.0000000

Survived	0	0.0000000
Pclass	0	0.0000000
Name	0	0.0000000
Sex	0	0.0000000
SibSp	0	0.0000000
Parch	0	0.0000000
Ticket	0	0.0000000
Fare	0	0.0000000

Проведем тест Литтла.

```
1 | mcar_test(titanic)
```

A tibble: 1 × 4

statistic	df	p.value	missing.patterns
<dbl>	<dbl>	<dbl>	<int>
598.5422	43	0	5

В данном случае вероятность (p-value) наблюдать такие пропуски при условии, что нулевая гипотеза верна, вообще равна нулю, и у нас появляется ещё больше оснований отвергнуть предположение о полностью случайных пропусках.

Функция сообщает, что в данных выявлено пять паттернов и собственно именно этими зависимостями между пропусками и наблюдаемыми данными мы и пользовались в многомерных способах заполнения пропусков.

## Временная сложность алгоритма

Откроем блокнот с кодом на Питоне

## Сравнение алгоритмов

Важность сравнения эффективности алгоритмов очевидна. Менее очевидным является способ их сравнения.

С одной стороны, мы могли бы измерять *фактическое время* работы алгоритма, однако такое измерение очень сильно зависело бы, в частности, от характеристик конкретного компьютера и в конечном счете не имело бы смысла.

Более удобно измерять работу алгоритма в *количестве операций*, необходимых для выполнения задачи.

## Линейный и бинарный поиск

Вернёмся к рассмотрению алгоритмов линейного и бинарного поиска. Вначале приведем соответствующие функции.

```
1 def linear(arr, x):
2
3     # объявим счетчик количества операций
4     counter = 0
5
6     for i in range(len(arr)):
7
8         # с каждой итерацией будем увеличивать счетчик на единицу
9         counter += 1
10
11        if arr[i] == x:
12            return i, counter
```

```
1 def binary(arr, x):
2
3     # объявим счетчик количества операций
4     counter = 0
5
6     low, high = 0, len(arr) - 1
7
8     while low <= high:
9
10        # увеличиваем счетчик с каждой итерацией цикла
11        counter += 1
12
13        mid = low + (high - low) // 2
14
15        if arr[mid] == x:
16            return mid, counter
17
18        elif arr[mid] < x:
19            low = mid + 1
20
21        else:
22            high = mid - 1
23
24    return -1
```

Эти алгоритмы аналогичны рассмотренным ранее, за исключением того, что в данном случае мы также будем считать количество операций. Под операций в контексте поиска мы будем понимать сравнение искомого числа с очередным элементом массива.

Пусть даны два отсортированных массива из восьми и шестнадцати чисел.

```
1 arr8 = np.array([3, 4, 7, 11, 13, 21, 23, 28])
2 arr16 = np.array([3, 4, 7, 11, 13, 21, 23, 28, 29, 30, 31, 33, 36, 37, 39, 42])
3
4 len(arr8), len(arr16)
```

```
1 (8, 16)
```

Найдем в этих массивах индекс чисел 28 и 42 соответственно. Алгоритм линейного поиска ожидаемо справится за восемь и шестнадцать операций.

```
1 # первым результатом функции будет индекс искомого числа,
2 # вторым - количество операций сравнения
3 linear(arr8, 28), linear(arr16, 42)
```

```
1 ((7, 8), (15, 16))
```

Заметим, что это **худший случай** из возможных (worst case scenario), искомые числа

оказались на самом неудачном месте.

Алгоритму же бинарного поиска потребуется всего лишь четыре и пять операций соответственно.

```
1 | binary(arr8, 28), binary(arr16, 42)
```

```
1 | ((7, 4), (15, 5))
```

Из приведенного выше кода следует, что алгоритмы не только различаются по количеству затрачиваемых операций, но и рост количества вычислений с ростом объема данных происходит не одинаково.

Количество операций ( $k$ ) линейного поиска растет пропорционально количеству данных ( $n$ ), т.е.  $k = n$ , это *линейная зависимость*, в бинарном же поиске *зависимость логарифмическая*  $k = \log(n)$ .

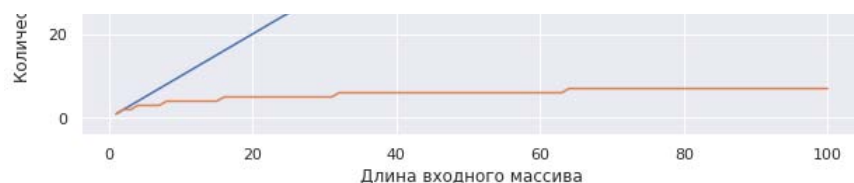
Сравним рост количества операций по мере роста объема данных.

```
1 | # посчитаем количество операций для входных массивов разной длины
2 | # создадим списки, куда будем записывать количество затраченных итераций
3 | ops_linear, ops_binary = [], []
4 |
5 | # будет 100 входных массивов длиной от 1 до 100 элементов
6 | input_arr = np.arange(1, 101)
7 |
8 | # на каждой итерации будем работать с массивом определенной длины
9 | for i in input_arr:
10 |
11 |     # внутри функций поиска создадим массив из текущего количества элементов
12 |     # и попросим найти последний элемент i - 1
13 |     _, l = linear(np.arange(i), i - 1)
14 |     _, b = binary(np.arange(i), i - 1)
15 |
16 |     # запишем количество затраченных операций в соответствующий список
17 |     ops_linear.append(l)
18 |     ops_binary.append(b)
```

Выведем результат на графике.

```
1 | plt.plot(input_arr, ops_linear, label = 'Линейный поиск')
2 | plt.plot(input_arr, ops_binary, label = 'Бинарный поиск')
3 |
4 | plt.title('Зависимость количества операций поиска от длины массива')
5 | plt.xlabel('Длина входного массива')
6 | plt.ylabel('Количество операций в худшем случае')
7 |
8 | plt.legend();
```





Вполне понятно, почему в алгоритме линейного поиска количество операций растет пропорционально объему данных. Выясним, почему количество операций бинарного поиска имеет логарифмическую зависимость.

## Логарифмическая зависимость бинарного поиска

На каждом этапе мы делим данные пополам и выполняем операцию сравнения. Например, массив из 16-ти чисел мы разделим четыре раза.

$$\frac{n}{2}, \frac{n}{4}, \frac{n}{8}, \frac{n}{16}, \dots, \frac{n}{2^k}$$

В нашем случае,

$$\frac{16}{2^k}$$

Выполнять такое деление мы будем до тех пор, пока количество чисел не будет равно одному (т.е. мы не найдем нужное число).

$$\frac{n}{2^k} = 1, 2^k = n$$

Тогда  $k$  (количество операций) будет равно

$$k = \log_2(n)$$

В нашем случае,

$$4 = \log_2(16)$$

## Нотация «О» большого

Такая оценка называется **временной сложностью** (time complexity) алгоритма и обычно выражается в **О-символике** или **нотации «О» большого** (big-O notation). Для приведенных выше алгоритмов запись будет выглядеть так:

- линейный поиск:  $O(n)$ ;
- бинарный поиск:  $O(\log(n))$ .

Повторюсь, что учитывается количество операций в худшем случае. При желании, можно сравнить и наиболее благоприятный сценарий (best case scenario). Тогда,

- для алгоритма линейного поиска его сложность в лучшем случае (если искомое значение стоит на первом месте) составит  $O(1)$ , то есть потребуется одна операция;

- при бинарном поиске, если искомое число находится, что наиболее удачно, в середине отсортированного массива, его сложность также составит  $O(1)$ .

Основание логарифма обычно не приводится поскольку

$$\log_a b = \frac{\log_c b}{\log_c a}$$

## Асимптотическое время

Нотация «О» большого отражает так называемое асимптотическое время работы алгоритма.

**Асимптотический анализ** (asymptotic analysis) изучает поведение функции при стремлении аргумента к определенному значению. В нашем случае мы смотрим на количество операций (поведение функции) при значительном увеличении объема данных (аргумент).

Скорость изменения функции ещё называют *порядком* изменения, от англ. order или нем. Ordnung, отсюда заглавная буква «О» в нотации.

Таким образом говорят, что в худшем случае линейный поиск соответствует линейному времени, а бинарный — логарифмическому.

Существует также константное время (когда сложность алгоритма не зависит от объема данных,  $O(1)$ ), квадратичное время  $O(n^2)$  или, например, факториальное  $O(n!)$ . В последнем случае, количество операций растет наиболее стремительно.

## Про константы

Внимательный читатель вероятно обратил внимание, что формула временной сложности бинарного поиска и фактическое количество операций приведенных алгоритмов не совпадают.

Мы взяли массивы из 8-ми и 16-ти чисел и по формуле, должны были уложиться в  $\log_2(8) = 3$  и  $\log_2(16) = 4$  операции. Фактически же мы затратили четыре и пять соответственно.

Другими словами сложность алгоритма бинарного поиска в худшем случае составляет  $O(\log(n) + 1)$ , однако так как нас интересует порядок (примерное понимание) роста функции при росте аргумента, а не точное значение, константы принято опускать.

Аналогично  $O(2 \log(n))$  сводится к  $O(\log(n))$ .

## Временная сложность в ML

Как применять вычислительную сложность алгоритма на практике? Рассмотрим два примера.

**Поиск ближайших соседей** по k-мерному дереву существенно быстрее, чем простой



перебор всех векторов сравнения и расчет расстояния:

- во втором случае (brute force) временная сложность равна  $O(n, m)$ , где  $n$  — количество векторов запроса, а  $m$  — количество векторов сравнения;
- в первом (kdtree) —  $O(n \log(m))$ , правда без учета операций для создания  $k$ -мерного дерева.

Кроме того, если мы знаем из-за какого компонента данных количество необходимых операций растет наиболее быстро, именно с этого компонента мы и начнем оптимизацию модели.

Например, из [документации](#) мы знаем, что **сложность алгоритма IterativeImputer** составляет

$$O(knp^3 \min(n, p))$$

где  $k$  — максимальное число итераций,  $n$  — количество наблюдений, а  $p$  — количество признаков. Таким образом, очевидно, что для снижения количества необходимых операций, в первую очередь нужно работать с признаками.

## Еще одно сравнение методов заполнения пропусков

### Способ сравнения

Выше мы сравнили способы заполнения пропусков, вначале используя одномерный метод или многомерную модель, а затем применив алгоритм логистической регрессии в задаче классификации.

Такой подход был выбран, чтобы показать, как работают методы заполнения пропусков в связке с реальной задачей машинного обучения. Для чистого же сравнения этих методов более правильным будет

- взять полный датасет (без пропущенных значений);
- случайным образом изъять часть значений;
- заполнить пропуски различными методами;
- сравнить результат с изъянными данными.

### Создание данных с пропусками

Возьмем уже знакомый нам по вводному курсу датасет, позволяющий классифицировать опухоли на доброкачественные и злокачественные.

```
1 # импортируем данные опухолей из модуля datasets библиотеки sklearn
2 from sklearn.datasets import load_breast_cancer
3
4 # выведем признаки и целевую переменную и поместим их в X_full и _ соответственно
5 X_full, _ = load_breast_cancer(return_X_y = True, as_frame = True)
```

```

6 |
7 | # масштабируем данные
8 | X_full = pd.DataFrame(StandardScaler().fit_transform(X_full), columns = X_full.columns)

```

Теперь напишем функцию, которая будет случайным образом добавлять пропуски в выбранные нами признаки.

```

1 | # нам понадобится модуль random
2 | import random
3 |
4 | # на вход функция будет получать полный датафрейм, номера столбцов признаков,
5 | # долю пропусков в каждом из столбцов и точку отсчета
6 | def add_nan(x_full, features, nan_share = 0.2, random_state = None):
7 |
8 |     random.seed(random_state)
9 |
10 |    # сделаем копию датафрейма
11 |    x_nan = x_full.copy()
12 |
13 |    # вначале запишем количество наблюдений и количество признаков
14 |    n_samples, n_features = x_full.shape
15 |
16 |    # посчитаем количество признаков в абсолютном выражении
17 |    how_many = int(nan_share * n_samples)
18 |
19 |    # в цикле пройдемся по номерам столбцов
20 |    for f in range(n_features):
21 |        # если столбец был указан в параметре features,
22 |        if f in features:
23 |            # случайным образом отберем необходимое количество индексов наблюдений (how_many
24 |            # из перечня, длиной с индекс (range(n_samples))
25 |            mask = random.sample(range(n_samples), how_many)
26 |            # заменим соответствующие значения столбца пропусками
27 |            x_nan.iloc[mask, f] = np.nan
28 |
29 |    # выведем датафрейм с пропусками
30 |    return X_nan

```

Обратите внимание на один нюанс. Функция **random.sample()**, что удобно, выбирает элемент *без возвращения*, то есть один раз выбрав наблюдение с индексом «один», второй раз она это наблюдение не выберет.

```

1 | # выведем пять чисел от 0 до 9
2 | random.seed(42)
3 | # с функцией random.sample() повторов не будет
4 | random.sample(range(10), 5)

```

```

1 | [1, 0, 4, 9, 6]

```

Мы могли бы использовать функции **np.random.randint()** или **random.choice()**, однако в этом случае из-за повторов процент пропусков был бы чуть ниже желаемого, например не 20, а 18.

```

1 | np.random.seed(42)
2 | # выберем случайным образом пять чисел от 0 до 9
3 | np.random.randint(0, 10, 5)

```

```

1 | array([6, 3, 7, 4, 6])

```

```

1 | random.seed(42)
2 | # выберем пять случайных чисел от 0 до 9
3 | [random.choice(range(10)) for _ in range(5)]

```

```

1 | [1, 0, 4, 3, 3]

```

Применим объявленную функцию к датафрейму и создадим 20 процентов пропусков в первом

столбце.

```
1 X_nan = add_nan(X_full,  
2                 features = [0],  
3                 nan_share = 0.2,  
4                 random_state = 42)
```

Проверим результат.

```
1 (X_nan.isna().sum() / len(X_nan)).round(2)
```

```
1 mean radius          0.2  
2 mean texture         0.0  
3 mean perimeter       0.0  
4 mean area            0.0  
5 mean smoothness      0.0  
6 mean compactness     0.0  
7 mean concavity       0.0  
8 mean concave points  0.0  
9 mean symmetry        0.0  
10 mean fractal dimension 0.0  
11 radius error        0.0  
12 texture error       0.0  
13 perimeter error     0.0  
14 area error          0.0  
15 smoothness error    0.0  
16 compactness error   0.0  
17 concavity error     0.0  
18 concave points error 0.0  
19 symmetry error      0.0  
20 fractal dimension error 0.0  
21 worst radius        0.0  
22 worst texture       0.0  
23 worst perimeter     0.0  
24 worst area          0.0  
25 worst smoothness    0.0  
26 worst compactness   0.0  
27 worst concavity     0.0  
28 worst concave points 0.0  
29 worst symmetry      0.0  
30 worst fractal dimension 0.0  
31 dtype: float64
```

Перейдем к заполнению пропусков.

## Заполнение пропусков

### Заполнение константой

```
1 # скопируем датасет  
2 fill_const = X_nan.copy()  
3 # заполним пропуски нулем  
4 fill_const.fillna(0, inplace = True)  
5 # убедимся, что пропусков не осталось  
6 fill_const.isnull().sum().sum()
```

```
1 '0'
```

### Заполнение медианой

```
1 # скопируем датасет
2 fill_median = X_nan.copy()
3 # заполним пропуски медианой
4 # по умолчанию, и .fillna(), и .median() работают со столбцами
5 fill_median.fillna(fill_median.median(), inplace = True)
6 # убедимся, что пропусков не осталось
7 fill_const.isnull().sum().sum()
```

```
1 '0'
```

Так как все признаки датасета количественные, не будем применять метод заполнения внутригрупповой медианой (для этого нам потребовалась бы хотя бы одна категориальная переменная).

## Заполнение линейной регрессией

Напишем функцию для заполнения пропусков линейной регрессией.

```
1 # передадим датафрейм, а также название столбца с пропусками
2 def linreg_imputer(df, col):
3
4     # обучающей выборкой будут строки без пропусков
5     train = df.dropna().copy()
6     # тестовой (или вернее выборкой для заполнения пропусков)
7     # будут те строки, в которых пропуски есть
8     test = df[df[col].isnull()].copy()
9
10    # выясним индекс столбца с пропусками
11    col_index = df.columns.get_loc(col)
12
13    # разделим "целевую переменную" и "признаки"
14    # обучающей выборки
15    y_train = train[col]
16    X_train = train.drop(col, axis = 1)
17
18    # из "тестовой" выборки удалим столбец с пропусками
19    test = test.drop(col, axis = 1)
20
21    # обучим модель линейной регрессии
22    model = LinearRegression()
23    model.fit(X_train, y_train)
24
25    # сделаем прогноз пропусков
26    y_pred = model.predict(test)
27    # вставим пропуски (value) на изначальное место (loc) столбца с пропусками (column)
28    test.insert(loc = col_index, column = col, value = y_pred)
29
30    # соединим датасеты и обновим индекс
31    df = pd.concat([train, test])
32    df.sort_index(inplace = True)
33
34    return df
```

Заполним пропуски.

```
1 fill_linreg = X_nan.copy()
2 fill_linreg = linreg_imputer(X_nan, 'mean radius')
3 fill_linreg.isnull().sum().sum()
```

```
1 '0'
```

## Заполнение с помощью MICE

```
1 fill_mice = X_nan.copy()
2 mice_imputer = IterativeImputer(initial_strategy = 'mean', # вначале заполним пропуски
3                                 estimator = LinearRegression(), # в качестве модели исг
4                                 random_state = 42 # добавим точку отсчета
5                                 )
6
7 # используем метод .fit_transform() для заполнения пропусков
8 fill_mice = pd.DataFrame(mice_imputer.fit_transform(fill_mice), columns = fill_mice.col
9 fill_linreg.isnull().sum().sum()
```

```
1 '0'
```

## Заполнение с помощью KNNImputer

```
1 fill_knn = X_nan.copy()
2
3 # используем те же параметры, что и раньше: пять "соседей" с одинаковыми весами
4 knn_imputer = KNNImputer(n_neighbors = 5, weights = 'uniform')
5
6 fill_knn = pd.DataFrame(knn_imputer.fit_transform(fill_knn), columns = fill_knn.columns
7 fill_knn.isnull().sum().sum()
```

```
1 '0'
```

## Оценка результата

Так как мы вычисляем количественные отклонения (разницу) прогнозного набора данных от полного, возведем отклонения в квадрат и суммируем вначале по столбцам (в нашем случае такой столбец один), а затем найдем общий квадрат всех отклонений датасета. Это и будет нашей метрикой.

Примечание. От RMSE метрика отличается только тем, что мы не делим на количество наблюдений (оно одинаковое в каждом датасете) и не извлекаем квадратный корень. Такое упрощение сделано для того, чтобы получившиеся в этом конкретном случае показатели было удобнее сравнивать.

```
1 # напишем функцию, которая считает сумму квадратов отклонений
2 # заполненного значения от исходного
3 def nan_mse(X_full, X_nan):
4     return ((X_full - X_nan)**2).sum().sum().round(2)
```

```
1 # создадим списки с датасетами и названиями методов
2 imputer = [fill_const, fill_median, fill_linreg, fill_mice, fill_knn]
3 name = ['constant', 'median', 'linreg', 'MICE', 'KNNImputer']
```

```
1 # в цикле оценим качество каждого из методов и выведем результат
2 for i, n in zip(imputer, name):
3     score = nan_mse(X_full, i)
4     print(n + ': ' + str(score))
```

```
1 constant: 122.7
2 median: 137.04
3 linreg: 0.03
4 MICE: 0.03
5 KNNImputer: 9.77
```

Используя более объективный метод сравнения, мы видим ещё бóльшую разницу результатов

одномерного и многомерного способов заполнения данных.

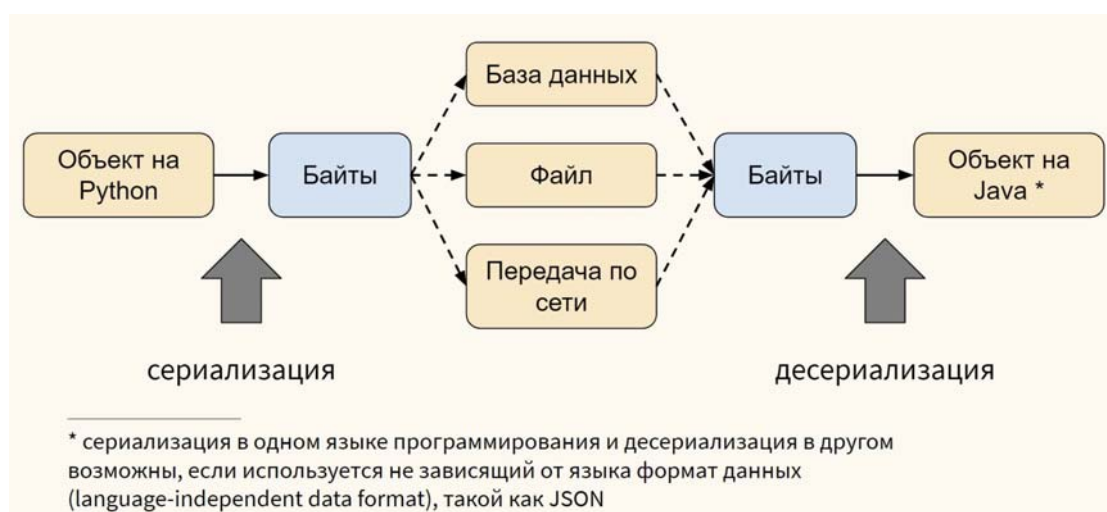
Обратите внимание, результат линейной регрессии совпал с результатом алгоритма MICE (как в случае с датасетом «Титаник», так и в текущем сравнении методов). Это логично, MICE, как и линейная регрессия использовали один и тот же алгоритм (класс `LinearRegression`) и одни и те же признаки без пропусков (и MICE не пришлось заполнять их средним значением).

## Сериализация и десериализация данных

Некоторые библиотеки создают файлы модели, в частности, в форматах JSON и Pickle. Зачем это нужно?

Благодаря этому у нас появляется возможность сохранить обученную модель на диске или передать ее по сети, а затем использовать для прогнозирования на новых данных *без необходимости повторного обучения*.

Процесс преобразования объекта (в частности, модели) в формат, пригодный для хранения и передачи данных, называется **сериализацией** (serialization). Процесс восстановления из этого формата называется соответственно **десериализацией** (deserialization).



Вначале рассмотрим в целом процесс сериализации в JSON и Pickle, а затем перейдем к работе с моделями.

## Формат JSON

Формат JSON расшифровывается как JavaScript Object Notation, но, несмотря на наличие слова JavaScript в своем названии, *не зависит от языка программирования*. Кроме того, он позволяет сериализовать данные в *человекочитаемый формат*.

Этот формат очень часто используется для передачи данных между клиентом и сервером. Например, запросим информацию о случайно выбранном банковском учреждении.

```
1 # импортируем модуль json,  
2 import json
```

```
3 # функцию urlopen() из модуля для работы с URL-адресами,
4 from urllib.request import urlopen
5 # а также функцию pprint() одноименной библиотеки
6 from pprint import pprint
7
8 url = 'https://random-data-api.com/api/v2/banks'
9
10 # получаем ответ (response) в формате JSON
11 with urlopen(url) as response:
12     # считываем его и закрываем объект response
13     data = response.read()
14
15 # данные пришли в виде последовательности байтов
16 print(type(data))
17 print()
18 # выполняем десериализацию
19 output = json.loads(data)
20 pprint(output)
21 print()
22 # и смотрим на получившийся формат
23 print(type(output))
```

```
1 <class 'bytes'>
2
3 {'account_number': '6695451080',
4  'bank_name': 'UBS CLEARING AND EXECUTION SERVICES LIMITED',
5  'iban': 'GB57WYQG90913422454081',
6  'id': 4063,
7  'routing_number': '014681402',
8  'swift_bic': 'BCYPGB2LXXX',
9  'uid': 'e1ef18ac-e286-471b-858f-f1a7cb3d85e6'}
10
11 <class 'dict'>
```

По сети мы получили объект типа bytes, то есть последовательность байтов, и десериализовали его с помощью **метода .loads()** в формат питоновского словаря.

Убедиться в человекочитаемости этого формата можно просто перейдя по ссылке в браузере: <https://random-data-api.com/api/v2/banks>.

## Вложенный словарь и список словарей

JSON-объект можно создать из вложенных питоновских словарей или списка словарей.

```
1 # создадим вложенные словари
2 sales = {
3     'PC' : {
4         'Lenovo' : 3,
5         'Apple' : 2
6     },
7     'Phone' : {
8         'Apple': 2,
9         'Samsung': 5
10    }
11 }
12
13 # и список из словарей
14 students = [
15     {
16         'id': 1,
17         'name': 'Alex',
18         'math': 5,
```

```
19         'computer science': 4
20     },
21     {
22         'id': 2,
23         'name': 'Mike',
24         'math': 4,
25         'computer science': 5
26     }
27 ]
```

## Методы .dumps() и .loads()

Вначале применим **метод .dumps()** для создания строкового JSON-объекта.

```
1 # преобразуем вложенный словарь в JSON
2 # дополнительно укажем отступ (indent)
3 json_sales = json.dumps(sales, indent = 4)
4
5 print(json_sales)
6 print(type(json_sales))
```

```
1 {
2     "PC": {
3         "Lenovo": 3,
4         "Apple": 2
5     },
6     "Phone": {
7         "Apple": 2,
8         "Samsung": 5
9     }
10 }
11 <class 'str'>
```

Обратите внимание, что хотя объект похож на питоновский словарь, тем не менее это строка.

Восстановим словарь с помощью **метода .loads()**.

```
1 sales = json.loads(json_sales)
2 print(sales)
3 print(type(sales))
```

```
1 {'PC': {'Lenovo': 3, 'Apple': 2}, 'Phone': {'Apple': 2, 'Samsung': 5}}
2 <class 'dict'>
```

## Методы .dump() и .load()

**Метод .dump()** создает последовательность байтов и используется для записи JSON-объекта в файл. Этот метод принимает два основных параметра: источник данных и файл, в который следует записать JSON-объект. Передадим ему этот файл с помощью конструкции with open().

```
1 # создадим файл students.json и откроем его для записи
2 with open('/content/students.json', 'w') as wf:
3     # поместим туда students, преобразовав в JSON
4     json.dump(students, wf, indent = 4)
```

Восстановим список словарей из файла.

```
1 # прочитаем файл из сессионного хранилища
2 with open('/content/students.json', "rb") as rf:
```



```
3 # и преобразуем обратно в список из словарей
4 students_out = json.load(rf)
```

Посмотрим на результат.

```
1 students_out
```

```
1 [{ 'id': 1, 'name': 'Alex', 'math': 5, 'computer science': 4},
2  { 'id': 2, 'name': 'Mike', 'math': 4, 'computer science': 5}]
```

Обратите внимание, результат десериализации — это новый объект.

```
1 print(students == students_out)
2 print(students is students_out)
```

```
1 True
2 False
```

## JSON и Pandas

Библиотека Pandas позволяет сохранить датафрейм в файл формата JSON, а также импортировать такой файл и соответственно восстановить датафрейм.

```
1 # импортируем датасет и преобразуем в датафрейм
2 from sklearn.datasets import load_breast_cancer
3 cancer, _ = load_breast_cancer(return_X_y = True, as_frame = True)
4
5 # создадим JSON-файл, поместим его в сессионное хранилище
6 cancer.to_json('/content/cancer.json')
7
8 # и сразу импортируем его и создадим датафрейм
9 pd.read_json('/content/cancer.json').head(3)
```

	mean radius	mean texture	mean perimeter	mean area	mean smoothness
0	17.99	10.38	122.8	1001.0	0.11840
1	20.57	17.77	132.9	1326.0	0.08474
2	19.69	21.25	130.0	1203.0	0.10960

Скриншот демонстрирует только часть датафрейма

## Pickle

Если JSON-объект можно создать на Питоне, а восстановить на Java, то объект Pickle сериализуется и десериализуется только с помощью Питона. За это отвечает одноименная библиотека.

```
1 import pickle
```

## Методы .dumps() и .loads()

Методы `.dumps()` и `.loads()` преобразуют объект в байты и восстанавливают исходный тип данных соответственно.

```
1 sales_pickle = pickle.dumps(sales)
2
3 print(sales_pickle)
4 print(type(sales_pickle))
```

```
1 b'\x80\x03}q\x00(X\x02\x00\x00\x00PCq\x01}q\x02(X\x06\x00\x00\x00Lenovoq\x03K\x03X\x05\
2 <class 'bytes'>
```

```
1 sales_out = pickle.loads(sales_pickle)
2
3 print(sales_out)
4 print(type(sales_out))
```

```
1 {'PC': {'Lenovo': 3, 'Apple': 2}, 'Phone': {'Apple': 2, 'Samsung': 5}}
2 <class 'dict'>
```

## Методы .dump() и .load()

Методы .dump() и .load() сериализуют объект в файл и соответственно десериализуют объект из файла.

```
1 # создадим файл students.p
2 # и откроем его для записи в бинарном формате (wb)
3 with open('/content/students.p', 'wb') as wf:
4     # поместим туда объект pickle
5     pickle.dump(students, wf)
```

```
1 # достанем этот файл из сессионного хранилища
2 # и откроем для чтения в бинарном формате (rb)
3 with open('/content/students.p', 'rb') as rf:
4     students_out = pickle.load(rf)
```

```
1 # выведем результат
2 students_out
```

```
1 [{'id': 1, 'name': 'Alex', 'math': 5, 'computer science': 4},
2  {'id': 2, 'name': 'Mike', 'math': 4, 'computer science': 5}]
```

При создании файла в формате pickle можно использовать расширения .p, .pkl или .pickle.

## Собственные объекты

В отличие JSON, который позволяет сериализовать только ограниченный набор объектов, Pickle может сериализовать любой питоновский объект, например собственную функцию или класс.

Начнем с **функций**.

```
1 # создадим функцию, которая будет выводить надпись "Some function!"
2 def foo():
3     print('Some function!')
4
5 # преобразуем эту функцию в объект Pickle
6 foo_pickle = pickle.dumps(foo)
7
8 # десериализуем и
9 foo_out = pickle.loads(foo_pickle)
10
11 # вызовем ее
12 foo_out()
```

```
1 | Some function!
```

То же самое можно сделать с **классом**.

```
1 | # создадим класс и объект этого класса
2 | class CatClass:
3 |
4 |     def __init__(self, color):
5 |         self.color = color
6 |         self.type_ = 'cat'
7 |
8 | Matroskin = CatClass('gray')
```

```
1 | # сериализуем класс в объект Pickle и поместим в файл
2 | with open('cat_instance.pkl', 'wb') as wf:
3 |     pickle.dump(Matroskin, wf)
```

```
1 | # достанем из файла и десериализуем
2 | with open('cat_instance.pkl', 'rb') as rf:
3 |     Matroskin_out = pickle.load(rf)
```

```
1 | # выведем атрибуты созданного нами объекта класса
2 | Matroskin_out.color, Matroskin_out.type_
```

```
1 | ('gray', 'cat')
```

Обратите внимание, что после десериализации мы смогли восстановить не только общую структуру класса, но и те данные (атрибут color), которые поместили в создаваемый объект.

Когда мы будем сохранять таким образом модели ML, вместо цвета шерсти животного мы будем хранить в Pickle, например, веса, полученные в результате обучения.

## Сохраняемость модели ML

Чаще всего для этого используется формат Pickle, потому что он способен сериализовать и десериализовать сложные объекты, которыми являются модели.

В машинном обучении говорят, что процесс сериализации и десериализации выполняется ради обеспечения **сохраняемости модели** (model persistence).

При этом, при использовании модуля Pickle, на этапах сериализации и десериализации важно работать с одинаковыми версиями Питона и используемых библиотек (например, sklearn). В противном случае результат может быть непредсказуемым. Это заставило некоторых специалистов предложить сериализацию модели ML с помощью JSON.

Рассмотрим сериализацию модели на практике. Вначале обучим модель логистической регрессии.

```
1 | # импортируем датасет о раке груди
2 | X, y = load_breast_cancer(return_X_y = True, as_frame = True)
3 |
4 | # импортируем класс для масштабирования данных,
5 | from sklearn.preprocessing import MinMaxScaler
6 | # функцию для разделения выборки на обучающую и тестовую части,
7 | from sklearn.model_selection import train_test_split
8 | # класс логистической регрессии
9 | from sklearn.linear_model import LogisticRegression
10 |
```

```
11 # разделим выборку
12 X_train, X_test, y_train, y_test = train_test_split(X, y,
13                                                    test_size = 0.30,
14                                                    random_state = 42)
15
16 # создадим объект класса MinMaxScaler
17 scaler = MinMaxScaler()
18
19 # масштабируем обучающую выборку
20 X_train_scaled = scaler.fit_transform(X_train)
21
22 # обучим модель на масштабированных train данных
23 model = LogisticRegression(random_state = 42).fit(X_train_scaled, y_train)
24
25 # используем минимальное и максимальное значения
26 # обучающей выборки для масштабирования тестовых данных
27 X_test_scaled = scaler.transform(X_test)
28
29 # сделаем прогноз
30 y_pred = model.predict(X_test_scaled)
```

Оценим результат.

```
1 # импортируем функцию для создания матрицы ошибок
2 from sklearn.metrics import confusion_matrix
3
4 # передадим матрице тестовые и прогнозные значения
5 # поменяем порядок так, чтобы злокачественные опухоли были положительным классом
6 model_matrix = confusion_matrix(y_test, y_pred, labels = [1,0])
7
8 # для удобства создадим датафрейм
9 model_matrix_df = pd.DataFrame(model_matrix)
10 model_matrix_df
```

	0	1
0	107	1
1	5	58

```
1 # рассчитаем accuracy
2 accuracy_score(y_test, y_pred).round(2)
```

```
1 0.96
```

Теперь выполним сериализацию и десериализацию модели.

```
1 # сериализуем и
2 with open('model.pickle', 'wb') as wf:
3     pickle.dump(model, wf)
```

```
1 # десериализуем модель
2 with open('model.pickle', 'rb') as rf:
3     model_out = pickle.load(rf)
```

Передадим десериализованной модели те же данные и убедимся, что она выдаст результат идентичный изначальной модели.

```
1 # напомним, десериализованная модель - это другой объект
2 y_pred_out = model_out.predict(X_test_scaled)
3
4 model_matrix = confusion_matrix(y_test, y_pred_out, labels = [1,0])
5 model_matrix_df = pd.DataFrame(model_matrix)
6 model_matrix_df
```

	0	1
0	107	1
1	5	58

```
1 | accuracy_score(y_test, y_pred).round(2)
```

```
1 | 0.96
```

## Подведем итог

В дополнительных материалах мы применили тест Литтла для выявления полностью случайных пропусков, познакомились с понятием временной сложности алгоритма, а также сериализацией и десериализацией модели ML.

---