

EDA данных о развитии малого и среднего предпринимательства Пиморского края

СОДЕРЖАНИЕ

[0 Импортирование библиотек](#)

[1 Просмотр и описание данных](#)

[1.1 Типы данных](#)

[1.2 Пропуски в данных](#)

[1.3 Дубликаты](#)

[1.4 Категориальные данные](#)

[1.5 Числовые данные](#)

[2 Тестирование гипотез](#)

[3 Итоговый вывод по результатам проведенного анализа](#)

0 ИМПОРТИРОВАНИЕ БИБЛИОТЕК

Библиотеки:

1. Библиотека `pandas` потребуется для работы с данными в табличном представлении.
2. Библиотека `pumpr` потребуется для создания данных в корреляционной матрице.
3. Библиотека `corpy` потребуется для создания глубоких копий объектов.
4. Библиотека `missingno` потребуется для визуализации пропусков в данных.
5. Библиотеки `matplotlib`, `seaborn` потребуются для визуализации данных.
6. Библиотека `scipy` потребуется для тестирования гипотез.
7. Библиотека `scikit_posthocs` потребуется для тестирования гипотез (тест Данна).

```
In [1]: import pandas as pd
import numpy as np
import copy
import missingno
import matplotlib.pyplot as plt
import seaborn as sns
from scipy.stats import normaltest
from scipy.stats import mstests
import scikit_posthocs as sp
```

Дополнительные настройки:

1. Настройка `pd.set_option('display.max_columns', None)` позволит вывести все колонки датафрейма не скрывая их.
2. Настройка `pd.set_option('display.float_format', lambda x: '%.2f' % x)` позволит вывести большие числа в числовом представлении а не в экспоненциальном.

```
In [2]: pd.set_option('display.max_columns', None)
pd.set_option('display.float_format', lambda x: '%.2f' % x)
```

1 ПРОСМОТР И ОПИСАНИЕ ДАННЫХ

В данном блоке будет осуществлен исследовательский анализ данных, описаны типы данных, пропуски, дубликаты, меры центральной тенденции а также взаимосвязи между переменными.

1.1 Типы данных

Подгружаем датасет.

```
In [3]: df = pd.read_csv('data_msp.csv') # Загрузка датасета
df.info() # Вывод информации по типам данных, количеству и пропускам
# Вывод первых 5 строк датасета

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 38366 entries, 0 to 38365
Data columns (total 53 columns):
 #   Column                Non-Null Count  Dtype  ---  --
 0   Наименование МСП      38366 non-null  object
 1   ИНН                   38366 non-null  object
 2   Муниципальное образование  38366 non-null  object
 3   Вид                   38366 non-null  object
 4   Категория             38366 non-null  object
 5   Код ОКВЭД             38366 non-null  object
 6   Вид деятельности      38366 non-null  object
 7   Доход                 33319 non-null  float64
 8   Расход               33319 non-null  float64
 9   ЕСХН                 21974 non-null  float64
10   УСН                  21974 non-null  float64
11   ЕНВД                 21974 non-null  float64
12   СРП                  21974 non-null  float64
13   Акции                7922 non-null   float64
14   Водный налог        7922 non-null   float64
15   Земельный налог      7922 non-null   float64
16   ЕСХН                 7922 non-null   float64
17   Задолженность и перерасчеты по отмененным налогам и 7922 non-null   float64
18   Земельный налог      7922 non-null   float64
19   Неналоговые доходы, администрируемые налоговыми органами и 7922 non-null   float64
20   НДС                 7922 non-null   float64
21   НДФЛ                7922 non-null   float64
22   ИФЛ                 7922 non-null   float64
23   Налог на имущество организаций и 7922 non-null   float64
24   Налог на прибыль     7922 non-null   float64
25   УСН                 7922 non-null   float64
26   Сборы за пользование объектами животного мира и за пользование объектами ВВР и 7922 non-null   float64
27   Страховые взносы мед и 7922 non-null   float64
28   Страховые взносы соц и 7922 non-null   float64
29   Страховые и другие взносы пенс и 7922 non-null   float64
30   Торговый сбор       7922 non-null   float64
31   Транспортный налог   7922 non-null   float64
32   Акции                35103 non-null  float64
33   Водный налог        35103 non-null  float64
34   ЕНВД                35103 non-null  float64
35   ЕСХН                35103 non-null  float64
36   Задолженность и перерасчеты по отмененным налогам у 35103 non-null  float64
37   Земельный налог     35103 non-null  float64
38   Неналоговые доходы, администрируемые налоговыми органами у 35103 non-null  float64
39   НДС                 35103 non-null  float64
40   ИФЛ                 35103 non-null  float64
41   НДФЛ               35103 non-null  float64
42   Налог на игорный бизнес у 35103 non-null  float64
43   Налог на имущество организаций у 35103 non-null  float64
44   Налог на прибыль     35103 non-null  float64
45   УСН                 35103 non-null  float64
46   Сборы за пользование объектами животного мира и за пользование объектами ВВР и 35103 non-null  float64
47   Страховые взносы мед и 35103 non-null  float64
48   Страховые взносы соц и 35103 non-null  float64
49   Страховые и другие взносы пенс у 35103 non-null  float64
50   Торговый сбор       35103 non-null  float64
51   Транспортный налог   35103 non-null  float64
52   Утилизационный сбор у 35103 non-null  float64
dtypes: float64(46), int64(1), object(6)
memory usage: 15.5+ MB
```

```
Out [3]:
```

	Наименование МСП	ИНН	Муниципальное образование	Вид	Категория	Код ОКВЭД	Вид деятельности	Доход	Расход	ЕС
0	ОБЩЕСТВО С ОГРАНИЧЕННОЙ ОТВЕТСТВЕННОСТЬЮ "ИПР"	701727301	ГОРОД ВЛАДИВОСТОК	Юридическое лицо	Микропредприятие	01.25.1	Выращивание прочих плодовых и яблочных культур	0.00	3000.00	0
1	ОБЩЕСТВО С ОГРАНИЧЕННОЙ ОТВЕТСТВЕННОСТЬЮ "ЧИП"	2537135020	ГОРОД ВЛАДИВОСТОК	Юридическое лицо	Микропредприятие	63.99.1	Деятельность в области консультирования и ин...	0.00	0.00	1
2	ОБЩЕСТВО С ОГРАНИЧЕННОЙ ОТВЕТСТВЕННОСТЬЮ "АИДАР"	2537141850	ГОРОД ВЛАДИВОСТОК	Юридическое лицо	Микропредприятие	46.15.4	Деятельность в области оптовой торговли радио...	0.00	0.00	0
3	ОБЩЕСТВО С ОГРАНИЧЕННОЙ ОТВЕТСТВЕННОСТЬЮ "АВТО..."	2537139280	ГОРОД ВЛАДИВОСТОК	Юридическое лицо	Микропредприятие	45.20	Техническое обслуживание и ремонт автотранспор...	нап	нап	1
4	ОБЩЕСТВО С ОГРАНИЧЕННОЙ ОТВЕТСТВЕННОСТЬЮ "САДО..."	2537138751	ГОРОД ВЛАДИВОСТОК	Юридическое лицо	Микропредприятие	71.12	Деятельность в области инженерных изысканий, и...	0.00	0.00	0

Датасет содержит 38366 записей с 53 столбцами, есть нулевые значения, пропуски, а также несоответствующие типы данных для переменных.

Приведем типы данных в соответствие:

1. Перенумеруем ИНН приведем к строковому значению так как это не число а идентификационный номер.
2. Для столбцов **Вид**, **Категория**, которые являются категориальными признаками, т.е. признаками, которые принимают ограниченное а обычно фиксированное количество возможных значений, установим тип данных как категория.

```
In [4]: # Приводим к типу данных строка
df['ИНН'] = df['ИНН'].astype('str')

# Приводим к типу данных категория
df['Категория'] = df['Категория'].astype('category')
df['Муниципальное образование'] = df['Муниципальное образование'].astype('category')
```

Проверим, прошли ли изменения.

```
In [5]: df.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 38366 entries, 0 to 38365
Data columns (total 53 columns):
 #   Column                Non-Null Count  Dtype  ---  --
 0   Наименование МСП      38366 non-null  object
 1   ИНН                   38366 non-null  object
 2   Муниципальное образование  38366 non-null  object
 3   Вид                   38366 non-null  object
 4   Категория             38366 non-null  object
 5   Код ОКВЭД             38366 non-null  object
 6   Вид деятельности      38366 non-null  object
 7   Доход                 33319 non-null  float64
 8   Расход               33319 non-null  float64
 9   ЕСХН                 21974 non-null  float64
10   УСН                  21974 non-null  float64
11   ЕНВД                 21974 non-null  float64
12   СРП                  21974 non-null  float64
13   Акции                7922 non-null   float64
14   Водный налог        7922 non-null   float64
15   Земельный налог      7922 non-null   float64
16   ЕСХН                 7922 non-null   float64
17   Задолженность и перерасчеты по отмененным налогам и 7922 non-null   float64
18   Земельный налог      7922 non-null   float64
19   Неналоговые доходы, администрируемые налоговыми органами и 7922 non-null   float64
20   НДС                 7922 non-null   float64
21   НДФЛ                7922 non-null   float64
22   ИФЛ                 7922 non-null   float64
23   Налог на имущество организаций и 7922 non-null   float64
24   Налог на прибыль     7922 non-null   float64
25   УСН                 7922 non-null   float64
26   Сборы за пользование объектами животного мира и за пользование объектами ВВР и 7922 non-null   float64
27   Страховые взносы мед и 7922 non-null   float64
28   Страховые взносы соц и 7922 non-null   float64
29   Страховые и другие взносы пенс и 7922 non-null   float64
30   Торговый сбор       7922 non-null   float64
31   Транспортный налог   7922 non-null   float64
32   Акции                35103 non-null  float64
33   Водный налог        35103 non-null  float64
34   ЕНВД                35103 non-null  float64
35   ЕСХН                35103 non-null  float64
36   Задолженность и перерасчеты по отмененным налогам у 35103 non-null  float64
37   Земельный налог     35103 non-null  float64
38   Неналоговые доходы, администрируемые налоговыми органами у 35103 non-null  float64
39   НДС                 35103 non-null  float64
40   ИФЛ                 35103 non-null  float64
41   НДФЛ               35103 non-null  float64
42   Налог на игорный бизнес у 35103 non-null  float64
43   Налог на имущество организаций у 35103 non-null  float64
44   Налог на прибыль     35103 non-null  float64
45   УСН                 35103 non-null  float64
46   Сборы за пользование объектами животного мира и за пользование объектами ВВР и 35103 non-null  float64
47   Страховые взносы мед и 35103 non-null  float64
48   Страховые взносы соц и 35103 non-null  float64
49   Страховые и другие взносы пенс у 35103 non-null  float64
50   Торговый сбор       35103 non-null  float64
51   Транспортный налог   35103 non-null  float64
52   Утилизационный сбор у 35103 non-null  float64
dtypes: category(2), float64(46), object(5)
memory usage: 15.0+ MB
```

Удалил столбец **вид** так как для последующего анализа он не представляет ценности.

```
In [6]: del df['Вид']
```

1.2 Пропуски в данных

Посчитаем процентное соотношение пропусков и заполненных значений для всех колонок.

```
In [7]: for col in df.columns:
    pct_missing = df[col].isnull().mean()
    print(f'{col} - {pct_missing:.1%}')

Наименование МСП - 0.0%
ИНН - 0.0%
Муниципальное образование - 0.0%
Категория - 0.0%
Код ОКВЭД - 0.0%
Вид деятельности - 0.0%
Доход - 13.2%
Расход - 13.2%
ЕСХН - 42.7%
УСН - 42.7%
ЕНВД - 42.7%
СРП - 42.7%
Акции и - 79.4%
Водный налог и - 79.4%
ЕНВД и - 79.4%
ЕСХН и - 79.4%
Задолженность и перерасчеты по отмененным налогам и - 79.4%
Земельный налог и - 79.4%
Земельные доходы, администрируемые налоговыми органами и - 79.4%
НДС и - 79.4%
НДФЛ и - 79.4%
НДФЛ и - 79.4%
Налог на имущество организаций и - 79.4%
Налог на прибыль и - 79.4%
УСН и - 79.4%
Сборы за пользование объектами животного мира и за пользование объектами ВВР и - 79.4%
Страховые взносы мед и - 79.4%
Страховые взносы соц и - 79.4%
Страховые и другие взносы пенс и - 79.4%
Торговый сбор и - 79.4%
Транспортный налог и - 79.4%
Акции у - 8.5%
Водный налог у - 8.5%
ЕНВД у - 8.5%
ЕСХН у - 8.5%
Задолженность и перерасчеты по отмененным налогам у - 8.5%
Земельный налог у - 8.5%
Неналоговые доходы, администрируемые налоговыми органами у - 8.5%
НДС у - 8.5%
НДФЛ у - 8.5%
НДФЛ у - 8.5%
Налог на игорный бизнес у - 8.5%
Налог на имущество организаций у - 8.5%
Налог на прибыль у - 8.5%
УСН у - 8.5%
Сборы за пользование объектами животного мира и за пользование объектами ВВР у - 8.5%
Страховые взносы мед у - 8.5%
Страховые взносы соц у - 8.5%
Страховые и другие взносы пенс у - 8.5%
Торговый сбор у - 8.5%
Транспортный налог у - 8.5%
Утилизационный сбор у - 8.5%
```

Наибольшее количество пропусков составляют колонки:

- 79.4%:
 - Акции и
 - Водный налог и
 - ЕНВД и
 - ЕСХН и
 - Задолженность и перерасчеты по отмененным налогам и
 - Земельный налог и
 - Неналоговые доходы, администрируемые налоговыми органами и
 - НДС и
 - НДФЛ и
 - НДФЛ и
 - Налог на имущество организаций и
 - Налог на прибыль и
 - УСН и
 - Сборы за пользование объектами животного мира и за пользование объектами ВВР и
 - Страховые взносы мед и
 - Страховые взносы соц и
 - Страховые и другие взносы пенс и
 - Торговый сбор и
 - Транспортный налог и
- 42.7%:
 - ЕСХН
 - УСН
 - ЕНВД
 - СРП

По остальным колонкам пропуски составляют небольшую долю в данных.

Визуализируем матрицу пропусков.

```
In [8]: print("Визуализация пропущенных значений со всеми пропусками.")
plt.figure(figsize=(15, 21))
missingno.matrix(df, fontsize = 14) # Передаем в качестве аргументов датафрейм, а также указываем размер шрифта
plt.show()
```

Визуализация пропущенных значений со всеми пропусками.

```
In [9]: print("Колонки с наибольшим количеством пропусков.")
df.iloc[:, :12].columns # Берем срез от датафрейма и выводим названия колонок

Колонки с наибольшим количеством пропусков.
```

```
Out [9]: Index(['Акции и', 'Водный налог и', 'ЕНВД и', 'ЕСХН и',
              'Задолженность и перерасчеты по отмененным налогам и',
              'Земельный налог и',
              'Неналоговые доходы, администрируемые налоговыми органами и',
              'Налог на имущество организаций и',
              'Налог на прибыль и', 'УСН и', 'Сборы за пользование объектами животного мира и за пользование объектами ВВР и',
              'Страховые взносы мед и', 'Страховые взносы соц и',
              'Страховые и другие взносы пенс и', 'Торговый сбор и',
              'Транспортный налог и'],
             dtype='object')
```

```
In [10]: print("Остатки следующие колонки с наибольшим количеством пропусков.")
df["Доход"].data[df["Акции и"].notnull() & df["Доход"].notnull()], fontsize = 14) # Передаем в качестве аргументов датафрейм
plt.show() # а также указываем размер шрифта

Визуализация пропущенных значений, в случае если убрать наибольшую часть пропусков в датасете.
```

```
In [11]: print("Остатки следующие колонки с наибольшим количеством пропусков.")
df["Доход"].notnull().iloc[:, 6:12].columns # Берем срез от датафрейма и выводим названия колонок

Остатки следующие колонки с наибольшим количеством пропусков.
```

```
Out [11]: Index(['Доход', 'Расход', 'ЕСХН', 'УСН', 'ЕНВД', 'СРП'], dtype='object')
```

В дальнейшем, колонки с системами налогообложения ценности особой не представляют, а вот информация о доходе и расходе у субъектов малого и среднего предпринимательства может потребоваться. Почистим в них пропуски и сформируем итоговый датафрейм.

```
In [12]: print("Визуализация пропущенных значений, в случае если оставить только необходимую для дальнейшего анализа информацию.")
df["Доход"].notnull().iloc[:, 6:12].columns # Передаем в качестве аргументов датафрейм, а также указываем размер шрифта
plt.show()
```

Визуализация пропущенных значений, в случае если оставить только необходимую для дальнейшего анализа информацию.

```
In [13]: print("Колонки с наибольшим количеством пропусков.")
df[df["Акции и"].notnull() & df["Доход"].notnull().iloc[:, 8:12].columns

Колонки с наибольшим количеством пропусков.
```

```
Out [13]: Index(['ЕСХН', 'УСН', 'ЕНВД', 'СРП'], dtype='object')
```

Так как проводить анализ по данным, содержащим подавляющее число пропусков целесообразно, создадим датафрейм с максимально возможной информацией без пропусков для последующего анализа.

```
In [14]: data_for_analysis = copy.deepcopy(df[df["Акции и"].notnull() & df["Доход"].notnull()])

Проверим количество муниципальных образований в данных. Важно что бы их было 34, так как в дальнейшем планируется проверить гипотезу где потребуются данные по каждому муниципалитету.
```

```
In [15]: if len(set(data_for_analysis["Муниципальное образование"])) == 34:
    print("Данные содержат информацию по каждому муниципальному образованию")
else:
    print("Данные содержат информацию не по каждому муниципальному образованию")

Данные содержат информацию по каждому муниципальному образованию
```

1.3 Дубликаты

1.3.1 Наименование МСП

```
In [17]: if df["Наименование МСП"].duplicated().any() == True:
    print("Данные содержат дубликаты.")
else:
    print("Данные не содержат дубликаты.")

Данные содержат дубликаты.
```

Данные о наименовании субъекта МСП содержат дубликаты по причине того, что названия могут быть одинаковые. В рамках данного анализа информация не представляет особого интереса.

1.3.2 ИНН

```
In [18]: if df["ИНН"].duplicated().any() == True:
    print("Данные содержат дубликаты.")
else:
    print("Данные не содержат дубликаты.")

Данные не содержат дубликаты.
```

Данные об ИНН субъекта МСП не содержат дубликаты по причине того, что ИНН это уникальный номер, он не может быть один у двух и более юридических лиц или индивидуальных предпринимателей. По данной информации можно строить агрегаты, а также использовать в качестве уникального ключа, идентифицирующего конкретную запись о субъекте МСП.

1.3.3 Муниципальное образование

```
In [19]: if df["Муниципальное образование"].duplicated().any() == True:
    print("Данные содержат дубликаты.")
else:
    print("Данные не содержат дубликаты.")

Данные содержат дубликаты.
```

Данные о муниципальном образовании, в котором зарегистрирован субъект МСП содержат дубликаты по причине того, что муниципальное образование это территория, на которой может осуществлять свою деятельность более одного субъекта МСП. По данной категориальной переменной можно строить агрегированную информацию в разрезе муниципальных образований.

1.3.4 Категория

```
In [20]: if df["Категория"].duplicated().any() == True:
    print("Данные содержат дубликаты.")
else:
    print("Данные не содержат дубликаты.")

Данные содержат дубликаты.
```

Данные о категории, к которой отнесен субъект МСП содержат дубликаты по причине того, что к одной и той же категории может относиться более одного субъекта МСП. По данной категориальной переменной можно строить агрегированную информацию в разрезе категории МСП.

1.3.5 Код ОКВЭД

```
In [21]: if df["Код ОКВЭД"].duplicated().any() == True:
    print("Данные содержат дубликаты.")
else:
    print("Данные не содержат дубликаты.")

Данные содержат дубликаты.
```

Данные о коде ОКВЭД, указавном как основной вид деятельности у субъекта МСП, содержат дубликаты по причине того, что субъекты МСП могут осуществлять одну и теки виды деятельности. По данной категориальной переменной можно строить агрегированную информацию в разрезе ОКВЭД МСП (в данном анализе это не представляет особого интереса).

1.3.6 Вид деятельности

```
In [22]: if df["Вид деятельности"].duplicated().any() == True:
    print("Данные содержат дубликаты.")
else:
    print("Данные не содержат дубликаты.")

Данные содержат дубликаты.
```

Данные о наименовании вида деятельности у субъекта МСП содержат дубликаты по причине того, что субъекты МСП могут осуществлять одни и теки виды деятельности. По данной категориальной переменной можно строить агрегированную информацию в разрезе ОКВЭД МСП (в данном анализе это не представляет особого интереса).

1.4 Категориальные данные

Визуализируем распределение субъектов малого и среднего предпринимательства по количеству в разрезе муниципальных образований Приморского края.

```
In [23]: # Задаем размеры визуализации
plt.figure(figsize=(15, 21))

# Строим гистограмму где по оси ординат передаем наименование муниципального образования
sns.countplot(y="Муниципальное образование",
              data=data_for_analysis, # Передаем датафрейм
              order = data_for_analysis["Муниципальное образование"].value_counts().index) # Считаем количество
plt.show() # Выводим график

Визуализация распределения субъектов малого и среднего предпринимательства по количеству в разрезе муниципальных образований Приморского края.
```

```
In [24]: # Задаем размеры визуализации
plt.figure(figsize=(15, 21))

# Строим гистограмму где по оси ординат передаем категорию бизнеса
sns.countplot(y="Категория",
              data=data_for_analysis, # Передаем датафрейм
              order = data_for_analysis["Категория"].value_counts().index) # Считаем количество
plt.show() # Выводим график

Визуализация распределения субъектов малого и среднего предпринимательства по количеству в разрезе категории микропредприятий, затем идут малые предприятия и самую незначительную долю составляет средний бизнес.
```

1.5 Числовые данные

Выведем меры центральной тенденции.

2 ТЕСТИРОВАНИЕ ГИПОТЕЗ

В данном блоке ответим на следующие вопросы:

- Отличаются ли муниципальные образования по уровню налоговой платежеспособности у субъектов МСП?
- В каких муниципальных образованиях наиболее низкий уровень налоговой платежеспособности?

1. Отличаются ли муниципальные образования по уровню налоговой платежеспособности у субъектов МСП?

Н1 - Муниципальные образования Приморского края не отличаются по уровню налоговой платежеспособности субъектов МСП.
Н0 - Муниципальные образования Приморского края отличаются по уровню налоговой платежеспособности субъектов МСП.

Переведем все числовые данные в целочисленному формату для корректного проведения тестов.

```
In [39]: # На каждой итерации цикла берем колонку и приводим
# формат числа к целочисленному
for col_name in data_new_features.iloc[:, 3:]:
    data_new_features[col_name] = data_new_features[col_name].astype('int64')
```

```
In [40]: # Выводим типы данных и отображим первые 5 строк датафрейма
data_new_features.info()
data_new_features.head()

<class 'pandas.core.frame.DataFrame'>
Int64Index: 7397 entries, 0 to 38364
Data columns (total 11 columns):
 #   Column              Non-Null Count  Dtype
---  --
 0   ИНН                 7397 non-null   object
 1   Муниципальное образование  7397 non-null   category
 2   Категория           7397 non-null   category
 3   Доход                7397 non-null   int64
 4   Расход              7397 non-null   int64
 5   НалогИ догт         7397 non-null   int64
 6   Иные платежи догт   7397 non-null   int64
 7   НалогИ уплаты       7397 non-null   int64
 8   Иные платежи уплаты 7397 non-null   int64
 9   Страховые неуплата  7397 non-null   int64
10   Страховые уплаты    7397 non-null   int64
dtypes: category(2), int64(8), object(1)
memory usage: 594.0+ KB
```

```
Out [40]:
```

	ИНН	Муниципальное образование	Категория	Доход	Расход	НалогИ догт	Иные платежи догт	НалогИ уплаты	Иные платежи уплаты	Страховые неуплата	Страховые уплаты	
5	2537138078	ГОРОД ВЛАДИВОСТОК	Микропредприятие	10448000	10324000	91725	0	0	38810	0	0	12870
16	2537138507	ГОРОД ВЛАДИВОСТОК	Микропредприятие	1086000	7919000	0	0	0	0	67	0	147720
21	2537138352	ГОРОД ВЛАДИВОСТОК	Микропредприятие	17592000	17410000	0	0	58538	0	0	0	435171
31	2537135648	ГОРОД ВЛАДИВОСТОК	Микропредприятие	60251000	58506000	11	0	487455	0	273	0	363008
33	2537138053	ГОРОД ВЛАДИВОСТОК	Микропредприятие	0	0	640333	0	0	0	0	0	0

```
In [41]: # Создадим датафрейм для проведения теста
df_test_mo = copy.deepcopy(data_new_features[['Муниципальное образование', 'НалогИ догт']])
# Выведем первые 5 строк датафрейма
df_test_mo.head()
```

```
Out [41]:
```

	Муниципальное образование	НалогИ догт
5	ГОРОД ВЛАДИВОСТОК	91725
16	ГОРОД ВЛАДИВОСТОК	0
21	ГОРОД ВЛАДИВОСТОК	0
31	ГОРОД ВЛАДИВОСТОК	11
33	ГОРОД ВЛАДИВОСТОК	640333

Проверим распределение на нормальность. Используем тест нормальности Андерсона-Дарлинга, так как он работает с большими выборками.

```
In [42]: def normal_test(df, column_name, alpha):
    """
    Функция, проверяющая распределение на нормальность

    Аргументы функции:
    df - датафрейм
    column_name - наименование колонки для проверки нормальности распределения
    alpha - уровень значимости

    Возвращаемое значение:
    Функция возвращает статистику критерия и p-value,
    а также принятие или отклонение H0.

    """
    print('Тест нормальности Андерсона-Дарлинга:')

    # Проводим тест
    stat, pval = normaltest(df[column_name])

    # Выводим статистику критерия и p-value
    print('H-statistic:', f'{stat:.3f}')
    print('P-Value:', f'{pval:.10f}')

    # Проверка условия для принятия или отклонения H0
    if pval > alpha:
        print('Принимаем H0 - Данные распределены нормально.')
```

```
In [43]: normal_test(df_test_mo, 'НалогИ догт', 0.05)

Тест нормальности Андерсона-Дарлинга:
Statistic: 23925.015
P-Value: 0.0000000000
Отклоняем H0 - Данные распределены не нормально.
```

Признак 'НалогИ догт' распределен не нормально. Учитывая данный факт для сравнения выборок будем использовать непараметрический критерий Краскела-Уоллиса.

Разобьем данные на 34 группы (так как в Приморском крае 34 муниципальных образования) для тестирования и проведем тест Краскела-Уоллиса.

```
In [44]: def kruskallwallis_test(df, group_name, values, alpha):
    """
    Функция, реализующая проведение
    теста Краскела-Уоллиса

    Аргументы функции:
    df - датафрейм
    group_name - наименование колонки с значениями для группировки
    values - наименование колонки с непрерывной величиной
    alpha - уровень значимости

    Возвращаемое значение:
    Функция возвращает статистику критерия и p-value,
    а также принятие или отклонение H0.

    """
    # Создаем пустой словарь, куда поместим группы для тестирования
    groups = {}
    # На каждой итерации цикла берем уникальное название для последующей фильтрации
    for grp in df[group_name].unique():
        # Создаем ключ словаря присвоив ему название группы для тестирования
        # передав значения, полученные в результате фильтрации датафрейма
        # (фильтруем исходный датафрейм по названию ключа, оставив только соответствующие значения)
        groups[grp] = df[values][df[group_name] == grp].values

    # Создаем тестовые группы, оставив только значения словаря
    test_group = groups.values()

    print('Тест Краскела-Уоллиса:')
    # Проводим тест
    H, pval = stats.kruskallwallis(*test_group)

    # Выводим статистику критерия и p-value
    print('H-statistic:', f'{H:.3f}')
    print('P-Value:', f'{pval:.10f}')

    # Проверка условия для принятия или отклонения H0
    if pval > alpha:
        print('Принимаем H0 - Между группами нет значительных различий.')
```

```
In [45]: kruskallwallis_test(df_test_mo, 'Муниципальное образование', 'НалогИ догт', 0.05 )

Тест Краскела-Уоллиса:
H-statistic: 101.258
P-Value: 0.0000000074
Отклоняем H0 - Между группами существуют значительные различия.
```

По итогам проведения теста получены статистически значимые различия между медианными значениями (данный стат критерий проверяет равенство медианных значений в выборках, что так же является плюсом так как средние значения сильно завышены за счет естественных выбросов) по уплате налогов бизнесом в муниципальных образованиях Приморского края. Это подтверждается критически низким значением p-value. Вероятность совершения ошибки первого рода равна 0.0000000074.

Учитывая, что в данных присутствуют естественные выбросы они могли повлиять на результаты теста. Проверим данный факт. Уберем выбросы (оставим диапазон + 3 сигмы (не смотря на то, что данное правило применимо к данным, распределенным нормально, используем его для создания диапазона) и повторим тестирование.

```
In [46]: # Создаем среднее
m = df_test_mo.mean()
# Создаем стандартное отклонение
std = df_test_mo.std()
# Создаем нижнюю границу
lower_bound = int(round(m - (3 * std)))
# Создаем верхнюю границу
upper_bound = int(round(m + (3 * std)))
# Создаем датафрейм без выбросов
mo_remove_outliers = df_test_mo[df_test_mo['НалогИ догт'].between(lower_bound, # Указываем нижнюю границу
                                                                    upper_bound)] # Указываем верхнюю границу

print(f'Число записей асеро (len(df_test_mo)) записей (ей).')
print(f'Число записей без выбросов составляет (len(mo_remove_outliers)) записей (ей).')

Число записей без выбросов составляет 7397 записей (ей).
```

Сравним статистику после удаления выбросов.

```
In [47]: print(f'Среднее до: {df_test_mo["НалогИ догт"].mean()} / Среднее после: {mo_remove_outliers["НалогИ догт"].mean()}')
print(f'Медиана до: {df_test_mo["НалогИ догт"].median()} / Медиана после: {mo_remove_outliers["НалогИ догт"].median()}')
print(f'Мода до: {df_test_mo["НалогИ догт"].mode()[0]} / Мода после: {mo_remove_outliers["НалогИ догт"].mode()[0]}')
print(f'Стандартное отклонение до: {df_test_mo["НалогИ догт"].std()} / Стандартное отклонение после: {m_remove_outliers["НалогИ догт"].std()}')
print(f'Дисперсия до: {df_test_mo["НалогИ догт"].var()} / Дисперсия после: {mo_remove_outliers["НалогИ догт"].var()}')

Среднее до: 173309.4472083277 / Среднее после: 39963.469783197834
Медиана до: 18.0 / Медиана после: 17.5
Мода до: 0 / Мода после: 0
Стандартное отклонение до: 4141243.9524956928 / Стандартное отклонение после: 364253.812159629
Дисперсия до: 17149901474082.148 / Дисперсия после: 132680839672.82228
```

После удаления выбросов значительно изменилось среднее. Это обусловлено удалением из данных экстремально высоких значений касательно неуплаты налогов, соответственно снизился разброс данных (стандартное отклонение и дисперсия). Медиана почти не изменилась, мода осталась неизменной. Так как критерий Краскела-Уоллиса проверяет равенство медианных значений в выборках удаление выбросов сильно не сказалось на изменении мер центральной тенденции, на основании которых проводится расчет, поэтому можно провести повторный тест по датфрейму без выбросов.

Проверим полученное распределение на нормальность. </p>
</div>
<div data-bbox="23 418 983 432" data-label="Text">
<pre>In [48]: normal_test(mo_remove_outliers, 'НалогИ догт', 0.05)

Тест нормальности Андерсона-Дарлинга:
Statistic: 14693.609
P-Value: 0.0000000000
Отклоняем H0 - Данные распределены не нормально.</pre>
</div>
<div data-bbox="103 433 536 436" data-label="Text">
<p>Снова разобьем данные на 34 группы и проведем тестирование.</p>
</div>
<div data-bbox="103 437 399 441" data-label="Text">
<p>H0 - Группы не отличаются.

H1 - Группы имеют существенные отличия.</p>
</div>
<div data-bbox="23 442 983 463" data-label="Text">
<pre>In [49]: kruskallwallis_test(mo_remove_outliers, 'Муниципальное образование', 'НалогИ догт', 0.05)

Тест Краскела-Уоллиса:
H-statistic: 103.747
P-Value: 0.0000000031
Отклоняем H0 - Между группами существуют значительные различия.</pre>
</div>
<div data-bbox="103 464 983 478" data-label="Text">
<p>Проведя два теста по данным с выбросами и по данным, взятым в диапазоне + 3 сигма установлено, что при частичном удалении выбросов вероятность получить статистически значимые различия в группах возрастает, но учитывая низкую вероятности совершения ошибки первого рода проведя тестирование на данных с выбросами, отбросив данных от выбросов в данном анализе не повлияла на результат.</p>
</div>
<div data-bbox="103 479 199 482" data-label="Section-Header">
<h4>С выбросами:</h4>
</div>
<div data-bbox="103 483 239 497" data-label="Text">
<p>H-statistic: 101.258

P-Value: 0.0000000074

+ - 3 сигмы:

P-Value: 0.0000000031</p>
</div>
<div data-bbox="103 498 149 501" data-label="Section-Header">
<h4>Вывод:</h4>
</div>
<div data-bbox="103 502 983 516" data-label="Text">
<p>На основании проведенного теста можно сделать вывод, что все муниципальные образования Приморского края отличаются по уровню налоговой платежеспособности субъектов малого и среднего предпринимательства. Это подтверждается низким уровне p-value, полученным в ходе тестирования. Вероятность совершения ошибки первого рода равна 0.0000000074.</p>
</div>
<div data-bbox="103 517 826 520" data-label="Section-Header">
<h3>2. В каких муниципальных образованиях наиболее низкий уровень налоговой платежеспособности?</h3>
</div>
<div data-bbox="103 521 983 535" data-label="Text">
<p>Для того, чтобы ответить на данный вопрос, а также используя полученную информацию по результатам предыдущего теста (установлены различия между муниципалитетами по уровню неуплаты налогов бизнесом) объединим муниципалитеты в следующие группы:</p>

low - медианная неуплата налога на одного предпринимателя меньше или равна 500 рублей.
medium - медианная неуплата налога на одного предпринимателя больше 500 рублей и меньше или равна 1500 рублей.
high - медианная неуплата налога на одного предпринимателя больше 1500 рублей.

</div>
<div data-bbox="103 536 983 547" data-label="Text">
<p>Создадим датафрейм с наименованием муниципального образования и медианным значением неуплаты налога в данном муниципалитете.</p>
</div>
<div data-bbox="23 548 983 570" data-label="Text">
<pre>In [50]: # Группируем данные по муи образованиям, агрегируем данные по медиане
mo_median_pay = df_test_mo.groupby('Муниципальное образование').agg(['НалогИ догт': ['median', 'count'
]])
Формируем новый датафрейм объединив информацию по муи образованиям и медиане
mo_merge = mo_median_pay['Муниципальное образование'].reset_index().merge(mo_median_pay['НалогИ догт'],
.reset_index())
Удалим столбец с индексом
del mo_merge['index']
Приводим данные по медиане к целочисленному формату
mo_merge['median'] = mo_merge['median'].astype('int64')
Сортируем данные по медиане по убыванию
mo_merge.sort_values('median', ascending=False)</pre>
</div>
<div data-bbox="23 571 983 603" data-label="Text">
<pre>Out [50]:
<table>
<thead>
<tr>
<th>Муниципальное образование</th>
<th>median</th>
<th>count</th>
</tr>
</thead>
<tbody>
<tr>
<td>17</td>
<td>РАЙОН ЛАЗОВСКИЙ</td>
<td>8787</td>
<td>8</td>
</tr>
<tr>
<td>12</td>
<td>РАЙОН АНЧИНСКИЙ</td>
<td>3686</td>
<td>3</td>
</tr>
<tr>
<td>32</td>
<td>РАЙОН ШКОТОВСКИЙ</td>
<td>1510</td>
<td>16</td>
</tr>
<tr>
<td>9</td>
<td>ГОРОД СПАСКО-ДАЛЬНИЙ</td>
<td>1462</td>
<td>30</td>
</tr>
<tr>
<td>24</td>
<td>РАЙОН ПОЖАРСКИЙ</td>
<td>1000</td>
<td>31</td>
</tr>
<tr>
<td>14</td>
<td>РАЙОН КАВАЛЕВСКИЙ</td>
<td>943</td>
<td>20</td>
</tr>
<tr>
<td>22</td>
<td>РАЙОН ПАРТИЗАНСКИЙ</td>
<td>907</td>
<td>22</td>
</tr>
<tr>
<td>13</td>
<td>РАЙОН ДАЛЬНЕРЕЧЕНСКИЙ</td>
<td>710</td>
<td>8</td>
</tr>
<tr>
<td>25</td>
<td>РАЙОН СПАССКИЙ</td>
<td>588</td>
<td>9</td>
</tr>
<tr>
<td>23</td>
<td>РАЙОН ПОГРЯНЧНЫЙ</td>
<td>501</td>
<td>10</td>
</tr>
<tr>
<td>11</td>
<td>РАЙОН ОЛЫГИНСКИЙ</td>
<td>469</td>
<td>10</td>
</tr>
<tr>
<td>21</td>
<td>ГОРОД ФОКИНО</td>
<td>383</td>
<td>21</td>
</tr>
<tr>
<td>6</td>
<td>ГОРОД ЛЕСОЗАВОДСК</td>
<td>249</td>
<td>30</td>
</tr>
<tr>
<td>29</td>
<td>РАЙОН ХОРОЛЬСКИЙ</td>
<td>244</td>
<td>11</td>
</tr>
<tr>
<td>2</td>
<td>ГОРОД БОЛЬШОЙ КАМЕНЬ</td>
<td>139</td>
<td>34</td>
</tr>
<tr>
<td>18</td>
<td>РАЙОН МИХАЙЛОВСКИЙ</td>
<td>135</td>
<td>12</td>
</tr>
<tr>
<td>8</td>
<td>ГОРОД ПАРТИЗАНСК</td>
<td>131</td>
<td>27</td>
</tr>
<tr>
<td>15</td>
<td>РАЙОН КИРОВСКИЙ</td>
<td>106</td>
<td>7</td>
</tr>
<tr>
<td>7</td>
<td>ГОРОД НАХОДКА</td>
<td>84</td>
<td>346</td>
</tr>
<tr>
<td>5</td>
<td>ГОРОД ДАЛЬНЕРЕЧЕНСК</td>
<td>68</td>
<td>38</td>
</tr>
<tr>
<td>28</td>
<td>РАЙОН ХАСАНСКИЙ</td>
<td>36</td>
<td>53</td>
</tr>
<tr>
<td>1</td>
<td>ГОРОД АРТЕМ</td>
<td>26</td>
<td>326</td>
</tr>
<tr>
<td>10</td>
<td>ГОРОД УСУРИЙСК</td>
<td>20</td>
<td>300</td>
</tr>
<tr>
<td>31</td>
<td>РАЙОН ЧУГУЙВСКИЙ</td>
<td>16</td>
<td>12</td>
</tr>
<tr>
<td>3</td>
<td>ГОРОД ВЛАДИВОСТОК</td>
<td>15</td>
<td>5760</td>
</tr>
<tr>
<td>0</td>
<td>ГОРОД АРСЕНЬЕВ</td>
<td>7</td>
<td>28</td>
</tr>
<tr>
<td>19</td>
<td>РАЙОН НАДЕЖДИНСКИЙ</td>
<td>7</td>
<td>79</td>
</tr>
<tr>
<td>16</td>
<td>РАЙОН КРАСНОАРМЕЙСКИЙ</td>
<td>6</td>
<td>17</td>
</tr>
<tr>
<td>4</td>
<td>ГОРОД ДАЛЬНЕГОРСК</td>
<td>3</td>
<td>73</td>
</tr>
<tr>
<td>26</td>
<td>РАЙОН ТЕРНЕЙСКИЙ</td>
<td>1</td>
<td>13</td>
</tr>
<tr>
<td>27</td>
<td>РАЙОН ХАНКАЙСКИЙ</td>
<td>0</td>
<td>8</td>
</tr>
<tr>
<td>20</td>
<td>РАЙОН ОКТЯБРЬСКИЙ</td>
<td>0</td>
<td>20</td>
</tr>
<tr>
<td>30</td>
<td>РАЙОН ЧЕРНЫГОРСКИЙ</td>
<td>0</td>
<td>10</td>
</tr>
<tr>
<td>33</td>
<td>РАЙОН ЯКОВЛЕВСКИЙ</td>
<td>0</td>
<td>5</td>
</tr>
</tbody>
</table>
</div>
<div data-bbox="103 604 696 607" data-label="Text">
<p>Присвоим каждому муниципальному образованию принадлежность к одной из трех групп.</p>
</div>
<div data-bbox="23 608 983 619" data-label="Text">
<pre>In [51]: mo_merge.loc[mo_merge['median'] <= 500, 'Категория'] = 'low'
mo_merge.loc[mo_merge['median'] > 500] & (mo_merge['median'] <= 1500), 'Категория'] = 'medium'
mo_merge.loc[mo_merge['median'] > 1500, 'Категория'] = 'high'</pre>
</div>
<div data-bbox="23 620 983 631" data-label="Text">
<pre>In [52]: mo_merge.sort_values('median', ascending=False)</pre>
</div>
<div data-bbox="23 632 983 664" data-label="Text">
<pre>Out [52]:
<table>
<thead>
<tr>
<th>Муниципальное образование</th>
<th>median</th>
<th>count</th>
<th>Категория</th>
</tr>
</thead>
<tbody>
<tr>
<td>17</td>
<td>РАЙОН ЛАЗОВСКИЙ</td>
<td>8787</td>
<td>8</td>
<td>high</td>
</tr>
<tr>
<td>12</td>
<td>РАЙОН АНЧИНСКИЙ</td>
<td>3686</td>
<td>3</td>
<td>high</td>
</tr>
<tr>
<td>32</td>
<td>РАЙОН ШКОТОВСКИЙ</td>
<td>1510</td>
<td>16</td>
<td>high</td>
</tr>
<tr>
<td>9</td>
<td>ГОРОД СПАСКО-ДАЛЬНИЙ</td>
<td>1462</td>
<td>30</td>
<td>medium</td>
</tr>
<tr>
<td>24</td>
<td>РАЙОН ПОЖАРСКИЙ</td>
<td>1000</td>
<td>31</td>
<td>medium</td>
</tr>
<tr>
<td>14</td>
<td>РАЙОН КАВАЛЕВСКИЙ</td>
<td>943</td>
<td>20</td>
<td>medium</td>
</tr>
<tr>
<td>22</td>
<td>РАЙОН ПАРТИЗАНСКИЙ</td>
<td>907</td>
<td>22</td>
<td>medium</td>
</tr>
<tr>
<td>13</td>
<td>РАЙОН ДАЛЬНЕРЕЧЕНСКИЙ</td>
<td>710</td>
<td>8</td>
<td>medium</td>
</tr>
<tr>
<td>25</td>
<td>РАЙОН СПАССКИЙ</td>
<td>588</td>
<td>9</td>
<td>medium</td>
</tr>
<tr>
<td>23</td>
<td>РАЙОН ПОГРЯНЧНЫЙ</td>
<td>501</td>
<td>10</td>
<td>medium</td>
</tr>
<tr>
<td>21</td>
<td>РАЙОН ОЛЫГИНСКИЙ</td>
<td>469</td>
<td>10</td>
<td>low</td>
</tr>
<tr>
<td>11</td>
<td>ГОРОД ФОКИНО</td>
<td>383</td>
<td>21</td>
<td>low</td>
</tr>
<tr>
<td>6</td>
<td>ГОРОД ЛЕСОЗАВОДСК</td>
<td>249</td>
<td>30</td>
<td>low</td>
</tr>
<tr>
<td>29</td>
<td>РАЙОН ХОРОЛЬСКИЙ</td>
<td>244</td>
<td>11</td>
<td>low</td>
</tr>
<tr>
<td>2</td>
<td>ГОРОД БОЛЬШОЙ КАМЕНЬ</td>
<td>139</td>
<td>34</td>
<td>low</td>
</tr>
<tr>
<td>18</td>
<td>РАЙОН МИХАЙЛОВСКИЙ</td>
<td>135</td>
<td>12</td>
<td>low</td>
</tr>
<tr>
<td>8</td>
<td>ГОРОД ПАРТИЗАНСК</td>
<td>131</td>
<td>27</td>
<td>low</td>
</tr>
<tr>
<td>15</td>
<td>РАЙОН КИРОВСКИЙ</td>
<td>106</td>
<td>7</td>
<td>low</td>
</tr>
<tr>
<td>7</td>
<td>ГОРОД НАХОДКА</td>
<td>84</td>
<td>346</td>
<td>low</td>
</tr>
<tr>
<td>5</td>
<td>ГОРОД ДАЛЬНЕРЕЧЕНСК</td>
<td>68</td>
<td>38</td>
<td>low</td>
</tr>
<tr>
<td>28</td>
<td>РАЙОН ХАСАНСКИЙ</td>
<td>36</td>
<td>53</td>
<td>low</td>
</tr>
<tr>
<td>1</td>
<td>ГОРОД АРТЕМ</td>
<td>26</td>
<td>326</td>
<td>low</td>
</tr>
<tr>
<td>10</td>
<td>ГОРОД УСУРИЙСК</td>
<td>20</td>
<td>300</td>
<td>low</td>
</tr>
<tr>
<td>31</td>
<td>РАЙОН ЧУГУЙВСКИЙ</td>
<td>16</td>
<td>12</td>
<td>low</td>
</tr>
<tr>
<td>3</td>
<td>ГОРОД ВЛАДИВОСТОК</td>
<td>15</td>
<td>5760</td>
<td>low</td>
</tr>
<tr>
<td>0</td>
<td>ГОРОД АРСЕНЬЕВ</td>
<td>7</td>
<td>28</td>
<td>low</td>
</tr>
<tr>
<td>19</td>
<td>РАЙОН НАДЕЖДИНСКИЙ</td>
<td>7</td>
<td>79</td>
<td>low</td>
</tr>
<tr>
<td>16</td>
<td>РАЙОН КРАСНОАРМЕЙСКИЙ</td>
<td>6</td>
<td>17</td>
<td>low</td>
</tr>
<tr>
<td>4</td>
<td>ГОРОД ДАЛЬНЕГОРСК</td>
<td>3</td>
<td>73</td>
<td>low</td>
</tr>
<tr>
<td>26</td>
<td>РАЙОН ТЕРНЕЙСКИЙ</td>
<td>1</td>
<td>13</td>
<td>low</td>
</tr>
<tr>
<td>27</td>
<td>РАЙОН ХАНКАЙСКИЙ</td>
<td>0</td>
<td>8</td>
<td>low</td>
</tr>
<tr>
<td>20</td>
<td>РАЙОН ОКТЯБРЬСКИЙ</td>
<td>0</td>
<td>20</td>
<td>low</td>
</tr>
<tr>
<td>30</td>
<td>РАЙОН ЧЕРНЫГОРСКИЙ</td>
<td>0</td>
<td>10</td>
<td>low</td>
</tr>
<tr>
<td>33</td>
<td>РАЙОН ЯКОВЛЕВСКИЙ</td>
<td>0</td>
<td>5</td>
<td>low</td>
</tr>
</tbody>
</table>
</div>
<div data-bbox="103 665 319 668" data-label="Text">
<p>Посмотрим на распределения в группах.</p>
</div>
<div data-bbox="23 669 983 683" data-label="Text">
<pre>In [53]: plt.figure(figsize=(28, 7))
plt.boxplot(x='median', y='Категория', data=mo_merge)
plt.show()</pre>
</div>
<div data-bbox="103 684 983 698" data-label="Text">
<p>Визуально из построенных боксплотов видно отличие групп. Подтвердим это проведя тест. Сформируем нулевую и альтернативную гипотезы.</p>
<p>Разобьем данные на 3 группы и проведем тестирование.</p>
<p>H0 - Группы не отличаются.

H1 - Группы имеют существенные отличия.</p>
</div>
<div data-bbox="23 699 983 710" data-label="Text">
<pre>In [54]: kruskallwallis_test(mo_merge, 'Категория', 'median', 0.05)

Тест Краскела-Уоллиса:
H-statistic: 21.136
P-Value: 0.0000257215
Отклоняем H0 - Между группами существуют значительные различия.</pre>
</div>
<div data-bbox="103 711 149 714" data-label="Section-Header">
<h4>Вывод:</h4>
</div>
<div data-bbox="103 715 983 729" data-label="Text">
<p>На основании проведенного теста можно сделать вывод, что между группами low, medium и high существуют различия по уровню налоговой платежеспособности субъектов малого и среднего предпринимательства. Это подтверждается низким уровнем p-value, полученным в ходе тестирования. Вероятность совершения ошибки первого рода равна 0.0000257215.</p>
</div>
<div data-bbox="103 730 9