


```
In [30]: msp_tax.columns = ['ИНН', 'Алименты у', 'Водный налог у', 'ЕНВД у', 'ЕСХН у', 'Задолженность и перерасчеты по отмененным налогам у', 'Земельный налог у', 'Налоговые доходы, администрируемые налоговыми органами у', 'НДС у', 'НДФЛ у', 'Налог на игорный бизнес у', 'Налог на имущество организаций у', 'Налог на прибыль у', 'УСН у', 'Сборы за пользование объектами животного мира и за пользование объектами ВВР у', 'Страховые взносы мед у', 'Страховые взносы соц у', 'Страховые и другие взносы пенс у', 'Торговый сбор у', 'Транспортный налог у', 'Утилизационный сбор у']

msp_tax.head()
```

Out [30]:

	ИНН	Алименты у	Водный налог у	ЕНВД у	ЕСХН у	Задолженность и перерасчеты по отмененным налогам у	Земельный налог у	Налоговые доходы, администрируемые налоговыми органами у	НДС у	НДФЛ у	Налог на имущество организаций у
0	814069428	0.0	0.0	0.0	0.0	0.0	4553.00	0.0	0.0	0.0	0.0
1	1433560214	0.0	0.0	0.0	0.0	0.0	0.00	0.0	0.0	0.0	0.0
2	1701050844	0.0	0.0	0.0	0.0	0.0	0.00	0.0	6270866.5	0.0	4956.0
3	1831167536	0.0	0.0	0.0	0.0	0.0	221487.77	0.0	0.0	0.0	0.0
4	2308178618	0.0	0.0	0.0	0.0	0.0	91889.00	0.0	0.0	0.0	0.0

5 rows × 12 columns

Выведем количество записей в датафрейме.

```
In [31]: print(f'Данные содержат {len(msp_tax)} записей (ей).')
Данные содержат 35103 записей (ей).
```

6 ЗАГРУЗКА ДАННЫХ ПО МСП БЕЗ УКАЗАНИЯ ОКВЭД

Загрузим ранее сформированный CSV файл, содержащий сведения о субъектах малого и среднего предпринимательства Приморского края, у которых не указан основной вид деятельности.

```
In [32]: # Создаем переменную names_5 с названиями колонок
names_6 = ['Наименование МСП', 'ИНН', 'Признак']

# Создаем датафрейм
activity_list = pd.read_csv('activity_list.csv'), # Указываем название файла
ser='q', # Указываем разделитель
encoding='utf-8', # Указываем кодировку
names=names_6, # Передаем названия колонок
engine='python') # Выбираем движок (python, так как в нем реализовано
# указание регулярного выражения в качестве разделителя
)

# Выводим первые пять строк
activity_list.head()
```

Out [32]:

	Наименование МСП	ИНН	Признак
0	ОБЩЕСТВО С ОГРАНИЧЕННОЙ ОТВЕТСТВЕННОСТЬЮ "ПУБ...	2502029222	Оклад не указан
1	"ХРЮП" ОБЩЕСТВО С ОГРАНИЧЕННОЙ ОТВЕТСТВЕННОСТЬЮ	2536013862	Оклад не указан
2	ОБЩЕСТВО С ОГРАНИЧЕННОЙ ОТВЕТСТВЕННОСТЬЮ "ЦЕНТ...	2536005173	Оклад не указан
3	ОБЩЕСТВО С ОГРАНИЧЕННОЙ ОТВЕТСТВЕННОСТЬЮ "АВИОР"	2536023155	Оклад не указан
4	ОБЩЕСТВО С ОГРАНИЧЕННОЙ ОТВЕТСТВЕННОСТЬЮ "НАУЧ...	2536006174	Оклад не указан

Сформируем датафрейм для последующего объединения с реестром.

```
In [33]: activity_list = activity_list[['ИНН', 'Признак']]
activity_list.head()
```

Out [33]:

	ИНН	Признак
0	2502029222	Оклад не указан
1	2536013862	Оклад не указан
2	2536005173	Оклад не указан
3	2536023155	Оклад не указан
4	2536006174	Оклад не указан

Выведем количество записей в датафрейме.

```
In [34]: print(f'Данные содержат {len(activity_list)} записей (ей).')
Данные содержат 135 записей (ей).
```

7 ФОРМИРОВАНИЕ ИТОГОВОГО ДАТАСЕТА

Соберем все данные в единый датасет.

```
In [35]: # Создаем единый датафрейм
data = reestr.merge(company_revenue, how='left')\
        .merge(company_tax_regime, how='left')\
        .merge(msp_attears, how='left')\
        .merge(msp_tax, how='left')\
        .merge(activity_list, how='left')

# Выводим первые пять строк
data.head()
```

Out [35]:

	Наименование МСП	ИНН	Муниципальное образование	Вид	Категория	Код ОКВЭД	Вид деятельности	Доход	Раско
0	СТЕПАНОВА ИРИНА	644934682794	ГОРОД ФОКИНО	Индивидуальный предприниматель	Микропредприятие	85.21	Ремонт электронной бытовой техники	NaN	Na
1	КУРЮН ЕВГЕНИЙ	870100455199	ГОРОД УССУРИЙСК	Индивидуальный предприниматель	Микропредприятие	48.32	Деятельность летного такси и аэродвигателей	NaN	Na
2	ГОДНОВА СВЕТЛАНА	870600600221	ГОРОД ВЛАДИВОСТОК	Индивидуальный предприниматель	Микропредприятие	47.11	Торговая розничная торговля пищевыми пр...	NaN	Na
3	НАВРОСЬ ДМИТРИЙ	870600537308	ГОРОД ВЛАДИВОСТОК	Индивидуальный предприниматель	Микропредприятие	85.29	Ремонт прочих предметов личного потребления и ...	NaN	Na
4	ОБЩЕСТВО С ОГРАНИЧЕННОЙ ОТВЕТСТВЕННОСТЬЮ "АГРО...	7017227901	ГОРОД ВЛАДИВОСТОК	Юридическое лицо	Микропредприятие	01.25.1	Выращивание прочих плодовых и ягодных культур	0.0	3000.

5 rows × 10 columns

Проверим, все ли данные корректно попали в итоговый датафрейм.

```
In [36]: data.info()

<class 'pandas.core.frame.DataFrame'>
Int64Index: 84335 entries, 4 to 83991
Data columns (total 54 columns):
# Column
Dtype
---
0 Наименование МСП
object
1 ИНН
int64
2 Муниципальное образование
object
3 Вид
object
4 Категория
object
5 Код ОКВЭД
object
6 Вид деятельности
object
7 Доход
float64
8 Расход
float64
9 ЕСХН
float64
10 УСН
float64
11 ЕНВД
float64
12 СПИ
float64
13 Алименты и
float64
14 Водный налог и
float64
15 ЕНВД и
float64
16 ЕСХН и
float64
17 Задолженность и перерасчеты по отмененным налогам и
float64
18 Земельный налог и
float64
19 Неналоговые доходы, администрируемые налоговыми органами и
float64
20 НДС и
float64
21 НДФЛ и
float64
22 НДФЛ и
float64
23 Налог на имущество организаций и
float64
24 Налог на прибыль и
float64
25 УСН и
float64
26 Сборы за пользование объектами животного мира и за пользование объектами ВВР и
float64
27 Страховые взносы мед и
float64
28 Страховые взносы соц и
float64
29 Страховые и другие взносы пенс и
float64
30 Торговый сбор и
float64
31 Транспортный налог и
float64
32 Алименты у
float64
33 Водный налог у
float64
34 ЕНВД у
float64
35 ЕСХН у
float64
36 Задолженность и перерасчеты по отмененным налогам у
float64
37 Земельный налог у
float64
38 Неналоговые доходы, администрируемые налоговыми органами у
float64
39 НДС у
float64
40 НДФЛ у
float64
41 НДФЛ у
float64
42 Налог на игорный бизнес у
float64
43 Налог на имущество организаций у
float64
44 Налог на прибыль у
float64
45 УСН у
float64
46 Сборы за пользование объектами животного мира и за пользование объектами ВВР у
float64
47 Страховые взносы мед у
float64
48 Страховые взносы соц у
float64
49 Страховые и другие взносы пенс у
float64
50 Торговый сбор у
float64
51 Транспортный налог у
float64
52 Утилизационный сбор у
float64
53 Признак
object
dtypes: float64(46), int64(1), object(7)
memory usage: 35.4+ MB
```

Разделим датафрейм на два (создадим глубокие копии датафрейма). Один будет содержать информацию по юридическим лицам второй по индивидуальным предпринимателям.

```
In [37]: legal_entity = copy.deepcopy(data[data['Вид'] == 'Юридическое лицо'])
entrepreneur = copy.deepcopy(data[data['Вид'] == 'Индивидуальный предприниматель'])
```

Посмотрим еще раз на данные.

```
In [38]: legal_entity.info()

<class 'pandas.core.frame.DataFrame'>
Int64Index: 38166 entries, 4 to 83991
Data columns (total 54 columns):
# Column
Dtype
---
0 Наименование МСП
object
1 ИНН
int64
2 Муниципальное образование
object
3 Вид
object
4 Категория
object
5 Код ОКВЭД
object
6 Вид деятельности
object
7 Доход
float64
8 Расход
float64
9 ЕСХН
float64
10 УСН
float64
11 ЕНВД
float64
12 СПИ
float64
13 Алименты и
float64
14 Водный налог и
float64
15 ЕНВД и
float64
16 ЕСХН и
float64
17 Задолженность и перерасчеты по отмененным налогам и
float64
18 Земельный налог и
float64
19 Неналоговые доходы, администрируемые налоговыми органами и
float64
20 НДС и
float64
21 НДФЛ и
float64
22 НДФЛ и
float64
23 Налог на имущество организаций и
float64
24 Налог на прибыль и
float64
25 УСН и
float64
26 Сборы за пользование объектами животного мира и за пользование объектами ВВР и
float64
27 Страховые взносы мед у
float64
28 Страховые взносы соц и
float64
29 Страховые и другие взносы пенс и
float64
30 Торговый сбор у
float64
31 Транспортный налог и
float64
32 Алименты у
float64
33 Водный налог у
float64
34 ЕНВД у
float64
35 ЕСХН у
float64
36 Задолженность и перерасчеты по отмененным налогам у
float64
37 Земельный налог у
float64
38 Неналоговые доходы, администрируемые налоговыми органами у
float64
39 НДС у
float64
40 НДФЛ у
float64
41 НДФЛ у
float64
42 Налог на игорный бизнес у
float64
43 Налог на имущество организаций у
float64
44 Налог на прибыль у
float64
45 УСН у
float64
46 Сборы за пользование объектами животного мира и за пользование объектами ВВР и
float64
47 Страховые взносы мед у
float64
48 Страховые взносы соц у
float64
49 Страховые и другие взносы пенс у
float64
50 Торговый сбор у
float64
51 Транспортный налог у
float64
52 Утилизационный сбор у
float64
53 Признак
object
dtypes: float64(46), int64(1), object(7)
memory usage: 16.1+ MB
```

Так как в итоговом датафрейме информация по индивидуальным предпринимателям в большинстве отсутствует, будем проводить последующий анализ по данным о юридических лицах. Кроме того приведем информацию по видам деятельности в соответствии с субъектом малого и среднего предпринимательства, где основной вид экономической деятельности не указан проставим соответствующий признак.

```
In [40]: legal_entity.loc[legal_entity['Признак'] == 'Оклад не указан', 'Код ОКВЭД'] = 'Не указан'
legal_entity.loc[legal_entity['Признак'] == 'Оклад не указан', 'Вид деятельности'] = 'Не указан'
```

Большая колонка признаков не понадобится, удалим ее.

```
In [41]: del legal_entity['Признак']

Выведем итоговый датасет.
```

```
In [42]: legal_entity.to_csv('data_msp.csv', index=False, encoding='utf-8')
```