

Предварительная обработка данных

Импортируем необходимые библиотеки. на данном этапе достаточно библиотеки pandas

In [136]:

```
import pandas
```

Исследуем файлы для подготовки их к формированию целевой базы данных

Подключаем файл customers_data.xls для анализа

In [137]:

```
df_customers_data = pandas.read_excel("a/xls/customers_data.xls")  
df_customers_data.info()
```

```
<class 'pandas.DataFrame'>
```

RangeIndex: 4977 entries, 0 to 4976
Data columns (total 5 columns):
Column Non-Null Count Dtype

0 Customer ID 4977 non-null str
1 Customer Name 4977 non-null str
2 Customer Payment Terms 4977 non-null str
3 Address 4977 non-null str
4 Credit Limit 4977 non-null int64
dtypes: int64(1), str(4)
memory usage: 194.5 KB

Анализируем данные в датафрейме df_customers_data

In [138]:

df_customers_data

Out[138]:

	Customer ID	Customer Name	Customer Payment Terms	Address	Credit Limit
0	C0200769623-0	WAL-MAR corp	NAH4	55599 Katherine Harbors Suite 551\nWest Brenda...	50000
1	C02009808	BEN E	NAD1	5488 Michael	50000

	Customer ID	Customer Name	Customer Payment Terms	Address	Credit Limit
	28-1			Inlet\nElizabethport, MP 17624	
2	C0200792734-2	MDV/ trust	NAA8	708 Taylor Cape\nJohnstad, MT 34743	100000
3	C0140105686-3	SYSC Ilc	CA10	4113 Dana Ridges\nEast Clarencestad, IA 61466	100000
4	C0140106181-4	WAL-MAR foundation	NAH4	2759 Kimberly Villages\nThompsons side, OR 79370	100000
...
4972	C75386549198-4995	Hull Inc	NATM	05134 Katherine Springs Suite 090\nJonathanvie...	20000
4973	C28188290658-4996	Ramirez-Cain	NATH	34213 Jennings Land Suite 112\nMaxwellchester,...	1000000

	Customer ID	Customer Name	Customer Payment Terms	Address	Credit Limit
4974	C20595783 803-4997	Sims PLC	NATU	91232 Brown Valleys\nPort Austin, MD 81384	20000
4975	C84713226 861-4998	Ferrell-Thomas	NATH	PSC 6545, Box 4898\nAPO AE 95400	100000
4976	C32572229 109-4999	Alvarez-Floyd	NAC6	69237 Smith Passage\nJaredfort, AS 91937	100000

4977 rows × 5 columns

Вывод: В датафрейме содержатся ID покупателя, имя клиента, условия оплаты, адрес и кредитный лимит. Из данных, представленных в датафрейме, можно предположить, что в данном датафрейме содержится данные связанные с дебиторской задолженностью.

Подсчитываем пустые значения по всем столбцам датафрейма df_customers_data

In [139]:

```
df_customers_data.isna().sum()
```

Out[139]:

```
Customer ID      0
Customer Name    0
Customer Payment Terms  0
Address          0
Credit Limit     0
dtype: int64
```

Подсчитываем количество дубликатов в датафрейме df_customers_data

In [140]:

```
df_customers_data.duplicated().sum()
```

Out[140]:

```
np.int64(0)
```

Исследуем выбросы

Метод describe() в библиотеке Pandas генерирует описательные статистики для столбцов DataFrame. Он вычисляет различные статистические показатели, такие как

количество записей, среднее значение, стандартное отклонение, минимум, максимум и перцентили для числовых столбцов. Также метод предоставляет обобщённые статистические данные для столбцов с типом данных объекта

In [141]:

```
df_customers_data.describe()
```

Out[141]:

Credit Limit	
count	4977.000000
mean	82194.092827
std	172857.225127
min	0.000000
25%	10000.000000
50%	20000.000000
75%	100000.000000
max	1000000.000000

In [142]:

```
df_customers_data['Credit Limit'].describe()
```

Out[142]:

```
count    4977.000000
mean     82194.092827
std      172857.225127
min        0.000000
25%      10000.000000
50%      20000.000000
75%     100000.000000
max     1000000.000000
Name: Credit Limit, dtype: float64
```

Вывод: Присутствует странный разброс значений в столбце 'Credit Limit', например минимальное нулевое значение. Кредитный лимит не может равняться нулю

Выведем нулевые значения

In [143]:

```
df_customers_data[df_customers_data['Credit Limit'] == 0]
```

Out[143]:

	Customer ID	Customer Name	Customer Payment Terms	Address	Credit Limit
94	C01000360 60-94	BIRD-I	NA10	1165 Robert Estate\nCar olshire, ME 37387	0
158	C02007935 68-158	NICH corp	NAA8	Unit 7163 Box 8239\nDPO AP 27060	0
253	C01000166 91-253	JAVA	NAA8	86645 Jennifer Mall\nMorg anfort, IN 12809	0
263	C01401044 75-263	FOOD trust	NAA8	69288 Cole Plains Suite 776\nNew Marcus, KS 65287	0
268	C01000057 65-268	HT HA foundation	NAA8	32322 Durham Hill\nPort Elizabethto wn, VA 14770	0
...
4179	C97244360 592-4179	PLAZA WA associates	NAG2	81540 Tammy Road Suite 296\nDuran bury, IL	0

	Customer ID	Customer Name	Customer Payment Terms	Address	Credit Limit
				94233	
4180	C72506571 732-4180	BARB associates	NAA8	196 Gary Street Apt. 917\nNew Joshuafort, PW 7...	0
4181	C37008718 163-4181	PLAZA WA llc	NAG2	206 Gregory Mountain\nMaloneport, OR 56931	0
4183	C67943642 515-4183	JALI	NAA8	106 Brown Curve Suite 797\nLawrenceville, ID 0...	0
4189	C19808856 070-4189	SOLOM associates	NAVE	890 Meredith Inlet Suite 110\nJennifer, H...	0

403 rows × 5 columns

Как видно в датафрейме присутствуют 403 строки содержащие ноль в столбце 'Credit Limit'. Удалим их

In [144]:

```
df_customers_data = df_customers_data[df_customers_data['Credit Limit'] != 0]
```

Проверяем

In [145]:

```
df_customers_data
```

Out[145]:

	Customer ID	Customer Name	Customer Payment Terms	Address	Credit Limit
0	C0200769623-0	WAL-MAR corp	NAH4	55599 Katherine Harbors Suite 551\nWest Brenda...	50000
1	C0200980828-1	BEN E	NAD1	5488 Michael Inlet\nEliza bethport, MP 17624	50000
2	C02007927	MDV/ trust	NAA8	708 Taylor Cape\nJoh	100000

	Customer ID	Customer Name	Customer Payment Terms	Address	Credit Limit
	34-2			nstad, MT 34743	
3	C01401056 86-3	SYSC Ilc	CA10	4113 Dana Ridges\nEa st Clarence sta d, IA 61466	100000
4	C01401061 81-4	WAL-MAR foundation	NAH4	2759 Kimberly Villages\nT hompsonsi de, OR 79370	100000
...
4972	C75386549 198-4995	Hull Inc	NATM	05134 Katherine Springs Suite 090\nJonat hanvie...	20000
4973	C28188290 658-4996	Ramirez-Ca in	NATH	34213 Jennings Land Suite 112\nMaxw ellchester,...	1000000
4974	C20595783 803-4997	Sims PLC	NATU	91232 Brown Valleys\nPo rt Austin, MD 81384	20000

	Customer ID	Customer Name	Customer Payment Terms	Address	Credit Limit
4975	C84713226 861-4998	Ferrell-Tho mas	NATH	PSC 6545, Box 4898\nAPO AE 95400	100000
4976	C32572229 109-4999	Alvarez-Flo yd	NAC6	69237 Smith Passage\nJ aredfort, AS 91937	100000

4574 rows × 5 columns

In [146]:

```
df_customers_data['Credit Limit'].describe()
```

Out[146]:

count	4574.000000
mean	89435.942282
std	178507.610874
min	5000.000000
25%	20000.000000
50%	50000.000000
75%	100000.000000
max	1000000.000000
Name: Credit Limit, dtype: float64	

Вывод: Как видно теперь минимальное значение по столбцу 'Credit Limit' составляет

5000 и общее количество записей в датафрейме составляет 4574

Подключаем файл payables_data.xls для анализа

In [147]:

```
df_payables_data = pandas.read_excel("a/xls/payables_data.xls")
df_payables_data.info()
```

```
<class 'pandas.DataFrame'>
RangeIndex: 20000 entries, 0 to 19999
Data columns (total 15 columns):
#   Column                                Non-Null Count  Dtype
---  -
0   Invoice Number                        20000 non-null  str
1   Posting Date                         20000 non-null  object
2   Invoice Date                         20000 non-null  object
3   Payment Date                        13988 non-null  object
4   Net Due Date (System Calculated Date) 20000 non-null  object
5   Supplier ID                         20000 non-null  str
6   Invoice Amount                       20000 non-null  int64
7   Fiscal year                         20000 non-null  str
8   Overdue                             20000 non-null  int64
9   Invoice Status                       20000 non-null  str
10  Spend Category                      20000 non-null  str
11  Total Outstanding amount             20000 non-null  int64
12  Late payment fees                   20000 non-null  int64
13  Payterm_n                           20000 non-null  int64
14  Vendor_Type                         20000 non-null  str
dtypes: int64(5), object(4), str(6)
memory usage: 2.3+ MB
```

Просматриваем данные

In [148]:

df_payables_data

Out[148]:

	In v o i c e N u m b e r	P o s t i n g D a t e	In v o i c e D a t e	P a y m e n t D a t e	N e t D u e D a t e (S y s t e m C a l c u l a t e d D a t e)	S u p p l i e r ID	In v o i c e A m o u n t	F i s c a l y e a r	O v e r d u e	In v o i c e S t a t u s	S p e n d C a t e g o r y	T o t a l O u t s t a n d i n g a m o u n t	L a t e p a y m e n t f e e s	P a y t e r m _ n	V e n d o r _ T y p e
0	IN V- 5	2 0 1	2 0 1	3 0- 1	2 0 2	S- 1 9	7 1 4	2 0 1	0	P a i d	Ta x e s	0	0	3 0	D o m

1

Invoice Number	Posting Date	Invoice Date	Payment Date	Net Due Date (System Calculated Date)	Supplier ID	Invoice Amount	Fiscal Year	Overdue	Invoice Status	Spend Category	Total Outstanding Amount	Late Payment fees	Payer - n	Vendor - Type
978675602067198	9-04-2020	9-04-2020	2-20-2020	0-04-2020	8	7	9-2020							estic
IN V-438	16-03-22	20-02-21	N a N	16-04-22	S-187	8575	2020-2	75	Unpai	R a w M at	8575	600	30	D o m es

Invoice Number	Posting Date	Invoice Date	Payment Date	Net Due Date (System Calculated Date)	Supplier ID	Invoice Amount	Fiscal Year	Overdue	Invoice Status	Spend Category	Total Outstanding Amount	Late Payment Fees	Payer - n	Vendor - Type
5639898658799	020	2-03-00:00:00		020			021		d	erial				tic
IN V-45868	26-12-2001	24-12-2001	20-12-2001	26-12-2002	S-300	6790	2019-2020	0	Paid	Raw Material	0	0	60	Domestic

Invoice Number	Posting Date	Invoice Date	Payment Date	Net Due Date (System Calculated Date)	Supplier ID	Invoice Amount	Fiscal Year	Overdue	Invoice Status	Spend Category	Total Outstanding Amount	Late Payment Fees	Payer - n	Vendor - Type
12911382721	9	9	1000000	0			0			ial				
INV-80916754	14-01-2020	20-01-2020	N a N	14-02-2020	S-497	6575	2020-2021	137	Unpaid	Services	6575	822	30	Domestic

Invoice Number	Posting Date	Invoice Date	Payment Date	Net Due Date (System Calculated Date)	Supplier ID	Invoice Amount	Fiscal Year	Overdue	Invoice Status	Spend Category	Total Outstanding Amount	Late Payment Fees	Payer - n	Vendor - Type
73911407		00:00												
INV-96237212873	2019-08-12:00	2019-08-12:00	2020-01-22:00	2020-08-03:00	S-310	12635	2019-2020	0	Paid	Raw Material	0	0	90	Domestic

19996

Invoice Number	Posting Date	Invoice Date	Payment Date	Net Due Date (System Calculated Date)	Supplier ID	Invoice Amount	Fiscal Year	Overdue	Invoice Status	Spend Category	Total Outstanding Amount	Late Payment fees	Payer - n	Vendor - Type
552310	00	00		00										
IN V-923263686839	2001-09-03 00:00	2001-09-03 00:00	2008-02-01 9	2001-09-03 00:00	S-449	13346	2019-2020	0	Paid	Raw Material	0	0	90	Domestic

Invoice Number	Posting Date	Invoice Date	Payment Date	Net Due Date (System Calculated Date)	Supplier ID	Invoice Amount	Fiscal Year	Overdue	Invoice Status	Spend Category	Total Outstanding amount	Late payment fees	Payment - n	Vendor - Type
893	0	0		0										
INV-2025-03-24	2025-03-24	2025-03-24	2025-03-24	2025-03-24	S-232	7437	2020-2021	113	Unpaid	Services	7437	904	30	Domestic

19998

5

IN
V-
5
3
0
6
0
6
0
2
5
9
6
0
4
0
9
2

2
0
1
9-
0
1-
0
5
0
0
0
0
0
0

2
0
1
9-
0
1-
0
5
0
0
0
0
0
0

2
7-
0
6-
2
0
1
9

Net Due Date (System Calculated Date)

S-
1
5
8

1
3
4
0
0

2
0
1
9-
2
0
2
0
0

0

P
a
i
d

R
a
w
M
a
t
e
r
i
a
l

0

0

9
0

D
o
m
e
s
t
i
c

In
v
o
i
c
e
N
u
m
b
e
r

P
o
s
t
i
n
g
D
a
t
e

In
v
o
i
c
e
D
a
t
e

P
a
y
m
e
n
t
D
a
t
e

S
u
p
p
l
i
e
r
I
D

In
v
o
i
c
e
A
m
o
u
n
t

F
i
s
c
a
l
y
e
a
r

O
v
e
r
d
u
e

In
v
o
i
c
e
S
t
a
t
u
s

S
p
e
n
d
C
a
t
e
g
o
r
y

T
o
t
a
l
O
u
t
s
t
a
n
d
i
n
g
a
m
o
u
n
t

L
a
t
e
p
a
y
m
e
n
t
f
e
e
s

P
a
y
t
e
r
m
_
n

V
e
n
d
o
r
_
T
y
p
e

IN V- 3 5 8 8 7 8 4 1 8 5 8 8 0 1 7 0	2 0 1 9- 0 1- 0 3 0 0: 0: 0 0 0	2 0 1 9- 0 1- 0 6 0 0: 0: 0 0 0	2 1 3- 0 4- 2 0 1 9	1 9- 0 1- 0 0: 0: 0 0 0	S- 1 6 9	1 0 3 2 0	2 0 1 9- 2 0 2 0	0	P a i d	R a w M a t e r i a l	0	0	9 0	D o m e s t i c
In v o i c e N u m b e r	P o s t i n g D a t e	In v o i c e D a t e	P a y m e n t D a t e	N e t D u e D a t e (S y s t e m C a l c u l a t e d D a t e)	S u p p l i e r I D	In v o i c e A m o u n t	F i s c a l y e a r	O v e r d u e	In v o i c e S t a t u s	S p e n d C a t e g o r y	T o t a l O u t s t a n d i n g a m o u n t	L a t e p a y m e n t f e e s	P a y t e r m - n	V e n d o r - T y p e

20000 rows × 15 columns

Вывод: В датафрейме содержатся данные о кредиторской задолженности клиентов.

Подсчитываем количество пустых значений в столбцах датафрейма

In [149]:

```
df_payables_data.isna().sum()
```

Out[149]:

Invoice Number	0	
Posting Date	0	
Invoice Date	0	
Payment Date	6012	
Net Due Date (System Calculated Date)	0	
Supplier ID	0	
Invoice Amount	0	
Fiscal year	0	
Overdue	0	
Invoice Status	0	
Spend Category	0	
Total Outstanding amount	0	
Late payment fees	0	
Payterm_n	0	
Vendor_Type	0	

dtype: int64

Вывод: Как видно из предыдущего блока, датафрейм df_payables_data содержит 6012 незаполненных данных (NaN) в столбце Payment Date, что свидетельствует о том, что в

данном столбце отсутствует дата последнего платежа клиентом.

Подсчитываем дубликаты

In [150]:

```
df_payables_data.duplicated().sum()
```

Out[150]:

```
np.int64(0)
```

Дублирование строк отсутствует

Как видно из анализа некоторые столбцы с датами содержат не верный формат. испрвим это. Преобразуем данные из формата object в datetime64

In [151]:

```
df_payables_data['Posting Date'] = pandas.to_datetime(df_payables_data['Posting Date'])
```

In [152]:

```
df_payables_data.info()
```

```
<class 'pandas.DataFrame'>
RangeIndex: 20000 entries, 0 to 19999
Data columns (total 15 columns):
#   Column                                Non-Null Count  Dtype
---  -
0   Invoice Number                        20000 non-null  str
1   Posting Date                         20000 non-null  datetime64[us]
2   Invoice Date                          20000 non-null  object
3   Payment Date                         13988 non-null  object
4   Net Due Date (System Calculated Date) 20000 non-null  object
5   Supplier ID                          20000 non-null  str
6   Invoice Amount                        20000 non-null  int64
7   Fiscal year                          20000 non-null  str
8   Overdue                              20000 non-null  int64
9   Invoice Status                        20000 non-null  str
10  Spend Category                       20000 non-null  str
11  Total Outstanding amount              20000 non-null  int64
12  Late payment fees                    20000 non-null  int64
13  Payterm_n                            20000 non-null  int64
14  Vendor_Type                          20000 non-null  str
dtypes: datetime64[us](1), int64(5), object(3), str(6)
memory usage: 2.3+ MB
```

In [153]:

```
df_payables_data['Invoice Date'] = pandas.to_datetime(df_payables_data['Invoice Date'])
df_payables_data.info()
```

```

<class 'pandas.DataFrame'>
RangeIndex: 20000 entries, 0 to 19999
Data columns (total 15 columns):
#   Column                                Non-Null Count  Dtype
---  -
0   Invoice Number                        20000 non-null  str
1   Posting Date                         20000 non-null  datetime64[us]
2   Invoice Date                         20000 non-null  datetime64[us]
3   Payment Date                        13988 non-null  object
4   Net Due Date (System Calculated Date) 20000 non-null  object
5   Supplier ID                         20000 non-null  str
6   Invoice Amount                       20000 non-null  int64
7   Fiscal year                         20000 non-null  str
8   Overdue                             20000 non-null  int64
9   Invoice Status                       20000 non-null  str
10  Spend Category                       20000 non-null  str
11  Total Outstanding amount             20000 non-null  int64
12  Late payment fees                    20000 non-null  int64
13  Payterm_n                           20000 non-null  int64
14  Vendor_Type                         20000 non-null  str
dtypes: datetime64[us](2), int64(5), object(2), str(6)
memory usage: 2.3+ MB

```

In [154]:

```

df_payables_data['Payment Date'] = pandas.to_datetime(df_payables_data['Payment Date'])
df_payables_data.info()

```

```

<class 'pandas.DataFrame'>
RangeIndex: 20000 entries, 0 to 19999
Data columns (total 15 columns):
#   Column                                Non-Null Count  Dtype
---  -
0   Invoice Number                        20000 non-null  str
1   Posting Date                         20000 non-null  datetime64[us]
2   Invoice Date                         20000 non-null  datetime64[us]
3   Payment Date                        13988 non-null  datetime64[us]
4   Net Due Date (System Calculated Date) 20000 non-null  object
5   Supplier ID                         20000 non-null  str

```

```

6 Invoice Amount          20000 non-null int64
7 Fiscal year            20000 non-null str
8 Overdue                20000 non-null int64
9 Invoice Status          20000 non-null str
10 Spend Category        20000 non-null str
11 Total Outstanding amount 20000 non-null int64
12 Late payment fees     20000 non-null int64
13 Payterm_n             20000 non-null int64
14 Vendor_Type           20000 non-null str
dtypes: datetime64[us](3), int64(5), object(1), str(6)
memory usage: 2.3+ MB

```

```

/var/folders/lq/kd04l4wx7tnd4251g2pdj9_40000gn/T/ipykernel_8311/1249894834.py:1:
UserWarning: Parsing dates in %d-%m-%Y format when dayfirst=False (the default) was
specified. Pass `dayfirst=True` or specify a format to silence this warning.
df_payables_data['Payment Date'] = pandas.to_datetime(df_payables_data['Payment Date'])

```

In [155]:

```

df_payables_data['Net Due Date (System Calculated Date)'] =
pandas.to_datetime(df_payables_data['Net Due Date (System Calculated Date)'])
df_payables_data.info()

```

```

<class 'pandas.DataFrame'>
RangeIndex: 20000 entries, 0 to 19999
Data columns (total 15 columns):
#   Column                                Non-Null Count  Dtype
---  -
0   Invoice Number                        20000 non-null  str
1   Posting Date                        20000 non-null  datetime64[us]
2   Invoice Date                        20000 non-null  datetime64[us]
3   Payment Date                       13988 non-null  datetime64[us]
4   Net Due Date (System Calculated Date) 20000 non-null  datetime64[us]
5   Supplier ID                        20000 non-null  str
6   Invoice Amount                      20000 non-null  int64
7   Fiscal year                        20000 non-null  str
8   Overdue                           20000 non-null  int64

```

```
9 Invoice Status          20000 non-null str
10 Spend Category        20000 non-null str
11 Total Outstanding amount 20000 non-null int64
12 Late payment fees      20000 non-null int64
13 Payterm_n             20000 non-null int64
14 Vendor_Type           20000 non-null str
dtypes: datetime64[us](4), int64(5), str(6)
memory usage: 2.3 MB
```

Проверяем данные

In [156]:

df_payables_data

Out[156]:

	In v o i c e N u m b e r	P o s t i n g D a t e	In v o i c e D a t e	P a y m e n t D a t e	N e t D u e D a t e (S y s t e m C a l c u l a t e d D a t e)	S u p p l i e r I D	In v o i c e A m o u n t	F i s c a l y e a r	O v e r d u e	In v o i c e S t a t u s	S p e n d C a t e g o r y	T o t a l O u t s t a n d i n g a m o u n t	L a t e p a y m e n t f e e s	P a y t e r m - n	V e n d o r - T y p e
0	IN V-59786756002067198	2019-04-12	2019-04-12	2019-02-03	2019-04-01	S-198	7147	2019-02-02	0	P a i d	T a x e s	0	0	30	D o m e s t i c
1	IN V-2020	2020	2020	N a	2020	S-1	85	2020	7	U n	R a	85	60	3	D o

INVOICE NUMBER	POSTING DATE	INVOICE DATE	PAYMENT DATE	NET DUE DATE (SYSTEM CALCULATED DATE)	SUPPLIER ID	INVOICE AMOUNT	FISCAL YEAR	OVERDUE	INVOICE STATUS	SPEND CATEGORY	TOTAL OUTSTANDING AMOUNT	LATE PAYMENT FEES	PAYER -	VENDOR - TYPE
438563988658658799	2020-03-01	2020-03-01	T	2020-04-01	87	75	2020-02-01	5	paid	W Material	75	0	0	mes tic
2	IN V-45	2020-03-01	2020-03-01	2020-03-01	S-30	679	2020-01-01	0	Paid	Raw M	0	0	60	Domes

Invoice Number	Posting Date	Invoice Date	Payment Date	Net Due Date (System Calculated Date)	Supplier ID	Invoice Amount	Fiscal Year	Overdue	Invoice Status	Spend Category	Total Outstanding Amount	Late Payment Fees	Payer - n	Vendor - Type
86812911382721	12-26-2016	12-24-2016	01-10-2017	02-26-2017	0	0	2020			Material				Domestic
INVOICE-80916	2020-01-1	2020-01-0	NAT	2020-02-1	S-497	6575	2020-2020	137	Unpaid	Services	6575	822	30	

In v o i c e N u m b e r	P o s t i n g D a t e	In v o i c e D a t e	P a y m e n t D a t e	N e t D u e D a t e (S y s t e m C a l c u l a t e d D a t e)	S u p p l i e r I D	In v o i c e A m o u n t	F i s c a l y e a r	O v e r d u e	In v o i c e S t a t u s	S p e n d C a t e g o r y	T o t a l O u t s t a n d i n g a m o u n t	L a t e p a y m e n t f e e s	P a y t e r m - n	V e n d o r - T y p e
75473911407	4	1		4			1							
IN V-96237212	2019-08-12	2019-04-12	2020-01-20	2020-02-03	S-310	12635	2019-02-20	0	P a i d	R a w M a t e r i a l	0	0	90	D o m e s t i c

In v o i c e N u m b e r	P o s t i n g D a t e	In v o i c e D a t e	P a y m e n t D a t e	N e t D u e D a t e (S y s t e m C a l c u l a t e d D a t e)	S u p p l i e r I D	In v o i c e A m o u n t	F i s c a l y e a r	O v e r d u e	In v o i c e S t a t u s	S p e n d C a t e g o r y	T o t a l O u t s t a n d i n g a m o u n t	L a t e p a y m e n t f e e s	P a y t e r m - n	V e n d o r - T y p e
87362728
19995	IN V-1713719	2019-05-12	2019-05-12	2020-05-02	S-206	14372	2019-2020	146	U n p a i d	R a w M a t e r i a l	14372	876	60	D o m e s t i c

19996

In v o i c e N u m b e r	P o s t i n g D a t e	In v o i c e D a t e	P a y m e n t D a t e	N e t D u e D a t e (S y s t e m C a l c u l a t e d D a t e)	S u p p l i e r I D	In v o i c e A m o u n t	F i s c a l y e a r	O v e r d u e	In v o i c e S t a t u s	S p e n d C a t e g o r y	T o t a l O u t s t a n d i n g a m o u n t	L a t e p a y m e n t f e e s	P a y t e r m _ n	V e n d o r _ T y p e
091552310	IN V-923263686	2001-09-03-06	2001-09-03-06	2001-09-03-09	S-449	13346	2019-2020	0	P a i d	R a w M a t e r i a l	0	0	90	D o m e s t i c

IN
V-
5
4
5
5
3
4
1
3
3
8
2
0

2020-09-02

Posting Date

2
0
2
0-
0
4-
0
2

Invoice Date

N
a
T

Payment Date	Amount	Balance
10/1/2023	100.00	100.00
10/15/2023	50.00	50.00
11/1/2023	75.00	25.00
11/15/2023	25.00	0.00
12/1/2023	0.00	0.00
12/15/2023	0.00	0.00
1/1/2024	0.00	0.00
1/15/2024	0.00	0.00
2/1/2024	0.00	0.00
2/15/2024	0.00	0.00
3/1/2024	0.00	0.00
3/15/2024	0.00	0.00
4/1/2024	0.00	0.00
4/15/2024	0.00	0.00
5/1/2024	0.00	0.00
5/15/2024	0.00	0.00
6/1/2024	0.00	0.00
6/15/2024	0.00	0.00
7/1/2024	0.00	0.00
7/15/2024	0.00	0.00
8/1/2024	0.00	0.00
8/15/2024	0.00	0.00
9/1/2024	0.00	0.00
9/15/2024	0.00	0.00
10/1/2024	0.00	0.00
10/15/2024	0.00	0.00
11/1/2024	0.00	0.00
11/15/2024	0.00	0.00
12/1/2024	0.00	0.00
12/15/2024	0.00	0.00
1/1/2025	0.00	0.00
1/15/2025	0.00	0.00
2/1/2025	0.00	0.00
2/15/2025	0.00	0.00
3/1/2025	0.00	0.00
3/15/2025	0.00	0.00
4/1/2025	0.00	0.00
4/15/2025	0.00	0.00
5/1/2025	0.00	0.00
5/15/2025	0.00	0.00
6/1/2025	0.00	0.00
6/15/2025	0.00	0.00
7/1/2025	0.00	0.00
7/15/2025	0.00	0.00
8/1/2025	0.00	0.00
8/15/2025	0.00	0.00
9/1/2025	0.00	0.00
9/15/2025	0.00	0.00
10/1/2025	0.00	0.00
10/15/2025	0.00	0.00
11/1/2025	0.00	0.00
11/15/2025	0.00	0.00
12/1/2025	0.00	0.00
12/15/2025	0.00	0.00
1/1/2026	0.00	0.00
1/15/2026	0.00	0.00
2/1/2026	0.00	0.00
2/15/2026	0.00	0.00
3/1/2026	0.00	0.00
3/15/2026	0.00	0.00
4/1/2026	0.00	0.00
4/15/2026	0.00	0.00
5/1/2026	0.00	0.00
5/15/2026	0.00	0.00
6/1/2026	0.00	0.00
6/15/2026	0.00	0.00
7/1/2026	0.00	0.00
7/15/2026	0.00	0.00
8/1/2026	0.00	0.00
8/15/2026	0.00	0.00
9/1/2026	0.00	0.00
9/15/2026	0.00	0.00
10/1/2026	0.00	0.00
10/15/2026	0.00	0.00
11/1/2026	0.00	0.00
11/15/2026	0.00	0.00
12/1/2026	0.00	0.00
12/15/2026	0.00	0.00
1/1/2027	0.00	0.00
1/15/2027	0.00	0.00
2/1/2027	0.00	0.00
2/15/2027	0.00	0.00
3/1/2027	0.00	0.00
3/15/2027	0.00	0.00
4/1/2027	0.00	0.00
4/15/2027	0.00	0.00
5/1/2027	0.00	0.00
5/15/2027	0.00	0.00
6/1/2027	0.00	0.00
6/15/2027	0.00	0.00
7/1/2027	0.00	0.00
7/15/2027	0.00	0.00
8/1/2027	0.00	0.00
8/15/2027	0.00	0.00
9/1/2027	0.00	0.00
9/15/2027	0.00	0.00
10/1/2027	0.00	0.00

2020-09-03

Net Due Date (System Calculated Date)

S-
2
3
2

Supplier ID

7
4
3
7

Invoice Amount

2
0
2
0-
2
0
2
1

Fiscal year

11
3

Overdue

Unpaid

Invoice Status

S
er
vi
ce
s

Spend Category

7
4
3
7

Total Outstanding amount

9
0
4

Late payment fees

30

Payterm –

Domestic

Vendor
Type

19998

In v o i c e N u m b e r	P o s t i n g D a t e	In v o i c e D a t e	P a y m e n t D a t e	N e t D u e D a t e (S y s t e m C a l c u l a t e d D a t e)	S u p p l i e r I D	In v o i c e A m o u n t	F i s c a l y e a r	O v e r d u e	In v o i c e S t a t u s	S p e n d C a t e g o r y	T o t a l O u t s t a n d i n g a m o u n t	L a t e p a y m e n t f e e s	P a y t e r m - n	V e n d o r - T y p e
5475														
IN V-53060602596040	2019-01-05	2019-01-05	2019-06-27	2019-01-08	S-158	13400	2019-2020	0	P a i d	R a w M a t e r i a l	0	0	90	D o m e s t i c

19999

92
IN
V-
3
5
8
8
7
8
4
1
8
5
8
8
0
1
7

2019-01-03

2019-02-05

2019-04-03

2019-01-06

S-169

10320

2019-02-02

0

Paid

Raw Material

0

0

90

Domestic


```
0 Invoice Number          20000 non-null str
1 Posting Date            20000 non-null datetime64[us]
2 Invoice Date             20000 non-null datetime64[us]
3 Payment Date            13988 non-null datetime64[us]
4 Net Due Date (System Calculated Date) 20000 non-null datetime64[us]
5 Supplier ID             20000 non-null str
6 Invoice Amount           20000 non-null int64
7 Fiscal year             20000 non-null str
8 Overdue                 20000 non-null int64
9 Invoice Status           20000 non-null str
10 Spend Category         20000 non-null str
11 Total Outstanding amount 20000 non-null int64
12 Late payment fees      20000 non-null int64
13 Payterm_n              20000 non-null int64
14 Vendor_Type            20000 non-null str
dtypes: datetime64[us](4), int64(5), str(6)
memory usage: 2.3 MB
```

Исследуем выбросы

In [158]:

```
df_payables_data.describe()
```

Out[158]:

	Posti ng Date	Invoi ce Date	Pay ment Date	Net Due Date (Syst em Calc ulate d Date)	Invoi ce A mo unt	Over due	Total Outs tandi ng amo unt	Late pay ment fees	Payt erm_ n
coun t	2000 0	2000 0	1398 8	2000 0	2000 0.000 000	2000 0.000 000	2000 0.000 000	2000 0.000 000	2000 0.000 000
mea n	2019- 09-12 17:50 :29.7 6000 0	2019- 09-09 14:59 :55.6 8000 0	2019- 08-29 08:18 :58.5 7592 2	2019- 11-16 04:20 :16.8 0000 0	1128 7.620 500	33.43 1150	3389. 9686 50	200.1 5435 0	61.00 5000
min	2019- 01-01 00:00 :00	2019- 01-01 00:00 :00	2019- 01-02 00:00 :00	2019- 01-02 00:00 :00	101.0 0000 0	0.000 000	0.000 000	0.000 000	0.000 000
25%	2019- 05-10 00:00 :00	2019- 05-04 00:00 :00	2019- 05-14 00:00 :00	2019- 07-02 00:00 :00	7295. 0000 00	0.000 000	0.000 000	0.000 000	30.00 0000
50%	2019- 09-09 12:00 :00	2019- 09-06 00:00 :00	2019- 08-22 00:00 :00	2019- 11-16 00:00 :00	1123 7.000 000	0.000 000	0.000 000	0.000 000	60.00 0000
75%	2019- 12-16 00:00 :00	2019- 12-14 00:00 :00	2019- 11-29 00:00 :00	2020- 03-07 00:00 :00	1562 4.000 000	48.00 0000	6044. 5000 00	246.0 0000 0	90.00 0000

	Posti ng Date	Invoi ce Date	Pay ment Date	Net Due Date (Syst em Calc ulate d Date)	Invoi ce Amo unt	Over due	Total Outs tandi ng amo unt	Late pay ment fees	Payt erm_ n
max	2020- 12-06 00:00 :00	2020- 12-06 00:00 :00	2020- 12-06 00:00 :00	2020- 12-09 00:00 :00	1999 9.000 000	266.0 0000 0	1998 1.000 000	2560. 0000 00	90.00 0000
std	NaN	NaN	NaN	NaN	5074. 4104 16	53.94 9846	5866. 8716 45	362.7 8822 2	24.68 9683

Вывод: Аномалии не выявлены

Анализируем файл receivables_data.xls.

In [159]:

```
df_receivables_data = pandas.read_excel("a/xls/receivables_data.xls")
df_receivables_data.info()
```

```

<class 'pandas.DataFrame'>
RangeIndex: 50000 entries, 0 to 49999
Data columns (total 19 columns):
#   Column                Non-Null Count  Dtype
---  -
0   Business Code          50000 non-null  str
1   Customer Number        50000 non-null  object
2   Customer Name          50000 non-null  str
3   Payment_Date           50000 non-null  object
4   Business Year          50000 non-null  int64
5   Posting_Date           50000 non-null  datetime64[us]
6   Due_Date               50000 non-null  datetime64[us]
7   Payterm                50000 non-null  int64
8   Invoice Currency        50000 non-null  str
9   Total Open Amount      50000 non-null  int64
10  USD_CURRENNCY          50000 non-null  str
11  Total Open Amount_USD  50000 non-null  float64
12  Customer Payment Terms 50000 non-null  str
13  Invoice ID              49994 non-null  float64
14  Is Open                50000 non-null  int64
15  DUNNLEVEL              50000 non-null  int64
16  Credit_limit           50000 non-null  int64
17  Baseline_Date          50000 non-null  datetime64[us]
18  Region                 50000 non-null  str
dtypes: datetime64[us](3), float64(2), int64(6), object(2), str(6)
memory usage: 7.2+ MB

```

In [160]:

```
df_receivables_data
```

Out[160]:

Region	Baseline Date	Credit Limit	DUNNLEVEL	IsOpen	InvoiceID	CustomerPaymentTerms	TotalOpenAmount_USD	USD_CURRENCY	TotalOpenAmount	InvoiceCurrency	Payterm	Due Date	Posting Date	Business Year	Payment Date	Customer Name	Customer Number	Business Code
WEST	2020-01-23	50000	0	0	1930438e+09	NAH4	54273.280	USD	54273	USD	15	2020-02-10	2020-02-11	2020	2020-11-00	WAL-MARcorp	200769623	U001
MIDWEST	2019-	5000	2	0	1929	NAD1	79656	USD	7965	USD	20	2019-	2019-	2019	2019-	BENE	20098	U001

Region	Baseline - Date	Credit - Limit	DUNNLEVEL	IsOpen	InvoiceID	CustomerPaymentTerms	TotalOpenAmount - USD	USD - CURRENCY	TotalOpenAmount	InvoiceCurrency	PayerTerm	Due - Date	Posting - Date	BusinessYear	Payment - Date	CustomerName	CustomerNumber	BusinessCode
ST	07-20	0			646e+09		.600		7			08-11	07-22		08-00:00:00		0828	
NORTHEAST	2019-09-1	100000	3	0	1.929874e+	NAA8	2253.860	USD	2254	USD	15	2019-09-2	2019-09-1	2019	2019-12-30	MDV / trust	200792734	U001

Region	Baseline - Date	Credit - Limit	DUNNLEVEL	Is Open	Invoice ID	Customer Payment Terms	Total Open Amount - USD	USD - CURRENCY	Total Open Amount	Invoice Currency	Pay term	Due - Date	Posting - Date	Business Year	Payment - Date	Customer Name	Customer Number	Business Code
SOUTHWEST	2020-03-26	100000	2	1	2960623e+09	CA10	2441.778	USD	3300	CAD	11	2020-04-10	2020-03-30	2020	01/00/1900	SYSC Ilc	140105686	CA02

[illegible]

Region	Baseline - Date	Credit - Limit	DUNN LEVEL	Is Open	Invoice ID	Customer Payment Terms	Total Open Amount - USD	USD - CURRENCY	Total Open Amount	Invoice Currency	Payer term	Due - Date	Posting - Date	Business Year	Payment - Date	Customer Name	Customer Number	Business Code
SOUTHEAST	2020-02-16	100000	1	0	1930537e+09	NA A 8	6120.860	USD	6121	USD	15	2020-03-05	2020-02-19	2020	2020-03-05	SAFEW associates	200772595	U001
					9							0:00:00						49997

[illegible]

Вывод: В файле представлены данные о дебиторской задолженности

Подсчитываем пустые значения в столбцах

In [161]:

```
df_receivables_data.isna().sum()
```

Out[161]:

Business Code	0
Customer Number	0
Customer Name	0
Payment_Date	0
Business Year	0
Posting_Date	0
Due_Date	0
Payterm	0
Invoice Currency	0
Total Open Amount	0
USD_CURRENNCY	0
Total Open Amount_USD	0
Customer Payment Terms	0
Invoice ID	6
Is Open	0
DUNNLEVEL	0
Credit_limit	0
Baseline_Date	0
Region	0

dtype: int64

Подсчитываем дубликаты

In [162]:

```
df_receivables_data.duplicated().sum()
```

Out[162]:

```
np.int64(36)
```

Есть дубликаты. Убираем их

In [163]:

```
df_receivables_data = df_receivables_data.drop_duplicates(keep=False)
```

Проверяем

In [164]:

```
df_receivables_data.duplicated().sum()
```

Out[164]:

np.int64(0)

In [165]:

df_receivables_data.info()

<class 'pandas.DataFrame'>

Index: 49928 entries, 0 to 49999

Data columns (total 19 columns):

#	Column	Non-Null Count	Dtype
0	Business Code	49928 non-null	str
1	Customer Number	49928 non-null	object
2	Customer Name	49928 non-null	str
3	Payment_Date	49928 non-null	object
4	Business Year	49928 non-null	int64
5	Posting_Date	49928 non-null	datetime64[us]
6	Due_Date	49928 non-null	datetime64[us]
7	Payterm	49928 non-null	int64
8	Invoice Currency	49928 non-null	str
9	Total Open Amount	49928 non-null	int64
10	USD_CURRENNCY	49928 non-null	str
11	Total Open Amount_USD	49928 non-null	float64
12	Customer Payment Terms	49928 non-null	str
13	Invoice ID	49922 non-null	float64
14	Is Open	49928 non-null	int64
15	DUNNLEVEL	49928 non-null	int64
16	Credit_limit	49928 non-null	int64
17	Baseline_Date	49928 non-null	datetime64[us]
18	Region	49928 non-null	str

dtypes: datetime64[us](3), float64(2), int64(6), object(2), str(6)

memory usage: 7.6+ MB

In [166]:

df_receivables_data

Out[166]:

										Region
										Baseline_Date
										Credit_Limit
										DUNNLEVEL
										IsOpen
										InvoiceID
										CustomerPaymentTerms
										TotalOpenAmount_USD
										USD_CURRENCY
										TotalOpenAmount
										InvoiceCurrency
										PayerTerm
										Due_Date
										Posting_Date
										BusinessYear
										Payment_Date
										CustomerName
										CustomerNumber
										BusinessCode
0	U0001	200769623	WAL-MAR Corp	2022-11-00	2020-01-26	2020-02-10	15	USD	54273	USD
0	U0001	200769623	WAL-MAR Corp	2022-11-00	2020-01-26	2020-02-10	15	USD	54273	USD
0	U0001	200769623	WAL-MAR Corp	2022-11-00	2020-01-26	2020-02-10	15	USD	54273	USD
0	U0001	200769623	WAL-MAR Corp	2022-11-00	2020-01-26	2020-02-10	15	USD	54273	USD
0	U0001	200769623	WAL-MAR Corp	2022-11-00	2020-01-26	2020-02-10	15	USD	54273	USD
0	U0001	200769623	WAL-MAR Corp	2022-11-00	2020-01-26	2020-02-10	15	USD	54273	USD
0	U0001	200769623	WAL-MAR Corp	2022-11-00	2020-01-26	2020-02-10	15	USD	54273	USD
0	U0001	200769623	WAL-MAR Corp	2022-11-00	2020-01-26	2020-02-10	15	USD	54273	USD
0	U0001	200769623	WAL-MAR Corp	2022-11-00	2020-01-26	2020-02-10	15	USD	54273	USD
0	U0001	200769623	WAL-MAR Corp	2022-11-00	2020-01-26	2020-02-10	15	USD	54273	USD
0	U0001	200769623	WAL-MAR Corp	2022-11-00	2020-01-26	2020-02-10	15	USD	54273	USD
0	U0001	200769623	WAL-MAR Corp	2022-11-00	2020-01-26	2020-02-10	15	USD	54273	USD
0	U0001	200769623	WAL-MAR Corp	2022-11-00	2020-01-26	2020-02-10	15	USD	54273	USD
0	U0001	200769623	WAL-MAR Corp	2022-11-00	2020-01-26	2020-02-10	15	USD	54273	USD
0	U0001	200769623	WAL-MAR Corp	2022-11-00	2020-01-26	2020-02-10	15	USD	54273	USD
0	U0001	200769623	WAL-MAR Corp	2022-11-00	2020-01-26	2020-02-10	15	USD	54273	USD
0	U0001	200769623	WAL-MAR Corp	2022-11-00	2020-01-26	2020-02-10	15	USD	54273	USD
0	U0001	200769623	WAL-MAR Corp	2022-11-00	2020-01-26	2020-02-10	15	USD	54273	USD
0	U0001	200769623	WAL-MAR Corp	2022-11-00	2020-01-26	2020-02-10	15	USD	54273	USD
0	U0001	200769623	WAL-MAR Corp	2022-11-00	2020-01-26	2020-02-10	15	USD	54273	USD
0	U0001	200769623	WAL-MAR Corp	2022-11-00	2020-01-26	2020-02-10	15	USD	54273	USD
0	U0001	200769623	WAL-MAR Corp	2022-11-00	2020-01-26	2020-02-10	15	USD	54273	USD
0	U0001	200769623	WAL-MAR Corp	2022-11-00	2020-01-26	2020-02-10	15	USD	54273	USD
0	U0001	200769623	WAL-MAR Corp	2022-11-00	2020-01-26	2020-02-10	15	USD	54273	USD
0	U0001	200769623	WAL-MAR Corp	2022-11-00	2020-01-26	2020-02-10	15	USD	54273	USD
0	U0001	200769623	WAL-MAR Corp	2022-11-00	2020-01-26	2020-02-10	15	USD	54273	USD
0	U0001	200769623	WAL-MAR Corp	2022-11-00	2020-01-26	2020-02-10	15	USD	54273	USD
0	U0001	200769623	WAL-MAR Corp	2022-11-00	2020-01-26	2020-02-10	15	USD	54273	USD
0	U0001	200769623	WAL-MAR Corp	2022-11-00	2020-01-26	2020-02-10	15	USD	54273	USD
0	U0001	200769623	WAL-MAR Corp	2022-11-00	2020-01-26	2020-02-10	15	USD	54273	USD
0	U0001	200769623	WAL-MAR Corp	2022-11-00	2020-01-26	2020-02-10	15	USD	54273	USD
0	U0001	200769623	WAL-MAR Corp	2022-11-00	2020-01-26	2020-02-10	15	USD	54273	USD
0	U0001	200769623	WAL-MAR Corp	2022-11-00	2020-01-26	2020-02-10	15	USD	54273	USD
0	U0001	200769623	WAL-MAR Corp	2022-11-00	2020-01-26	2020-02-10	15	USD	54273	USD
0	U0001	200769623	WAL-MAR Corp	2022-11-00	2020-01-26	2020-02-10	15	USD	54273	USD
0	U0001	200769623	WAL-MAR Corp	2022-11-00	2020-01-26	2020-02-10	15	USD	54273	USD
0	U0001	200769623	WAL-MAR Corp	2022-11-00	2020-01-26	2020-02-10	15	USD	54273	USD
0	U0001	200769623	WAL-MAR Corp	2022-11-00	2020-01-26	2020-02-10	15	USD	54273	USD
0	U0001	200769623	WAL-MAR Corp	2022-11-00	2020-01-26	2020-02-10	15	USD	54273	USD
0	U0001	200769623	WAL-MAR Corp	2022-11-00	2020-01-26	2020-02-10	15	USD	54273	USD
0	U0001	200769623	WAL-MAR Corp	2022-11-00	2020-01-26	2020-02-10	15	USD	54273	USD
0	U0001	200769623	WAL-MAR Corp	2022-11-00	2020-01-26	2020-02-10	15	USD	54273	USD
0	U0001	200769623	WAL-MAR Corp	2022-11-00	2020-01-26	2020-02-10	15	USD	54273	USD
0	U0001	200769623	WAL-MAR Corp	2022-11-00	2020-01-26	2020-02-10	15	USD	54273	USD
0	U0001	200769623	WAL-MAR Corp	2022-11-00	2020-01-26	2020-02-10	15	USD	54273	USD
0	U0001	200769623	WAL-MAR Corp	2022-11-00	2020-01-26	2020-02-10	15	USD	54273	USD
0	U0001	200769623	WAL-MAR Corp	2022-11-00	2020-01-26	2020-02-10	15	USD	54273	USD
0	U0001	200769623	WAL-MAR Corp	2022-11-00	2020-01-26	2020-02-10	15	USD	54273	USD
0	U0001	200769623	WAL-MAR Corp	2022-11-00	2020-01-26	2020-02-10	15	USD	54273	USD
0	U0001	200769623	WAL-MAR Corp	2022-11-00	2020-01-26	2020-02-10	15	USD	54273	USD
0	U0001	200769623	WAL-MAR Corp	2022-11-00	2020-01-26	2020-02-10	15	USD	54273	USD
0	U0001	200769623	WAL-MAR Corp	2022-11-00	2020-01-26	2020-02-10	15	USD	54273	USD
0	U0001	200769623	WAL-MAR Corp	2022-11-00	2020-01-26	2020-02-10	15	USD	54273	USD
0	U0001	200769623	WAL-MAR Corp	2022-11-00	2020-01-26	2020-02-10	15	USD	54273	USD
0	U0001	200769623	WAL-MAR Corp	2022-11-00	2020-01-26	2020-02-10	15	USD	54273	USD
0	U0001	200769623	WAL-MAR Corp	2022-11-00	2020-01-26	2020-02-10	15	USD	54273	USD
0	U0001	200769623	WAL-MAR Corp	2022-11-00	2020-01-26	2020-02-10	15	USD	54273	USD
0	U0001	200769623	WAL-MAR Corp	2022-11-00	2020-01-26	2020-02-10	15	USD	54273	USD
0	U0001	200769623	WAL-MAR Corp	2022-11-00	2020-01-26	2020-02-10	15	USD	54273	USD
0	U0001	200769623	WAL-MAR Corp	2022-11-00	2020-01-26	2020-02-10	15	USD	54273	USD
0	U0001	200769623	WAL-MAR Corp	2022-11-00	2020-01-26	2020-02-10	15	USD	54273	USD
0	U0001	200769623	WAL-MAR Corp	2022-11-00	2020-01-26	2020-02-10	15	USD	54273	USD
0	U0001	200769623	WAL-MAR Corp	2022-11-00	2020-01-26	2020-02-10	15	USD	54273	USD
0	U0001	200769623	WAL-MAR Corp	2022-11-00	2020-01-26	2020-02-10	15	USD	54273	USD
0	U0001	200769623	WAL-MAR Corp	2022-11-00	2020-01-26	2020-02-10	15	USD	54273	USD
0	U0001	200769623	WAL-MAR Corp	2022-11-00	2020-01-26	2020-02-10	15	USD	54273	USD
0	U0001	200769623	WAL-MAR Corp	2022-11-00	2020-01-26	2020-02-10	15	USD	54273	USD
0	U0001	200769623	WAL-MAR Corp	2022-11-00	2020-01-26	2020-02-10	15	USD	54273	USD
0	U0001	200769623	WAL-MAR Corp	2022-11-00	2020-01-26	2020-02-10	15	USD	54273	USD
0	U0001	200769623	WAL-MAR Corp	2022-11-00	2020-01-26	2020-02-10	15	USD	54273	USD
0	U0001	200769623	WAL-MAR Corp	2022-11-00	2020-01-26	2020-02-10	15	USD	54273	USD
0	U0001	200769623	WAL-MAR Corp	2022-11-00	2020-01-26	2020-02-10	15	USD	54273	USD
0	U0001	200769623	WAL-MAR Corp	2022-11-00	2020-01-26	2020-02-10	15	USD	54273	USD
0	U0001	200769623	WAL-MAR Corp	2022-11-00	2020-01-26	2020-02-10	15	USD	54273	USD
0	U0001	200769623	WAL-MAR Corp	2022-11-00	2020-01-26	2020-02-10	15	USD	54273	USD
0	U0001	200769623	WAL-MAR Corp	2022-11-00	2020-01-26	2020-02-10	15	USD	54273	USD
0	U0001	200769623	WAL-MAR Corp	2022-11-00	2020-01-26	2020-02-10	15	USD	54273	USD
0	U0001	200769623	WAL-MAR Corp	2022-11-00	2020					

Region	Baseline - Date	Credit - Limit	DUNN LEVEL	Is Open	Invoice ID	Customer Payment Terms	Total Open Amount - USD	USD - CURRENCY	Total Open Amount	Invoice Currency	Payer term	Due - Date	Posting - Date	Business Year	Payment - Date	Customer Name	Customer Number	Business Code
UTHWEST	20 - 03 - 26	0000			960623e + 09	10	41.778	D	00	D	1	20 - 04 - 10	20 - 03 - 30	20	/ 00 / 19 00	SC Ilc	0105686	02
WEST	2019 - 11 - 10	100000	3	0	1.930148e + 09	NAH4	33133.290	USD	33133	USD	15	2019 - 11 - 28	2019 - 11 - 13	2019	2019 - 11 - 25 00 :	WAL - MARR found a	200769623	U001

Customer Information										Payment Details										Invoice & Billing										Account & Location																				
Customer Data					Business Info					Payment Terms					Invoice Data					Billing Info					Account Details					Location																				
Customer Name		Customer Number			Business Name		Business Year			Payment Date		Posting Date			Business Year		Due Date			Pay Term		Invoice Currency			Total Open Amount		USD - CURRENCY			Total Open Amount - USD		Customer Payment Terms			Invoice ID		Is Open			DUNN LEVEL		Credit Limit			Baseline Date		Region			
CustName		CustNum			BusName		BusYear			PayDate		PostDate			BusYear		DueDate			PayTerm		InvCur			TotalOpen		USD_Cur			TotalOpen_USD		CustTerms			InvID		IsOpen			DUNNLEVEL		CreditLimit			BaselineDate		Region			
WAL-MAR Co		200769623			20019		2019			2019-08-15		2019-08-30			2019		2019-08-30			15		USD			6767		USD			6766.540		NAH4			1.929744e+09		0			2		1000000			2019-08-11		WEST			
SAFEW		20077			2020		2020			2020-		2020-			2020		2020-			15		USD			6121		USD			6120.		NAA8			1.930		0			1		100000			2020-		SOUTH			
U001		U001			U001		U001			U001		U001			U001		U001			U001		U001			U001		U001			U001			U001		U001			U001			U001		U001			U001		U001		
49996		49996			49996		49996			49996		49996			49996		49996			49996		49996			49996		49996			49996			49996		49996			49996			49996		49996			49996		49996		
4999		4999			4999		4999			4999		4999			4999		4999			4999		4999			4999		4999			4999			4999		4999			4999			4999		4999			4999		4999		

Region	Baseline - Date	Credit - Limit	DUNN LEVEL	Is Open	Invoice ID	Customer Payment Terms	Total Open Amount - USD	USD - CURRENCY	Total Open Amount	Invoice Currency	Payer term	Due - Date	Posting - Date	Business Year	Payment - Date	Customer Name	Customer Number	Business Code
EAST	02 - 16	0			537e + 09		860					03 - 05	02 - 19		03 - 05	associates	2595	7
MIDWEST	2019 - 11 - 2	100000	0	0	1930199e +	NA A8	63480	USD	63	USD	15	2019 - 12 - 1	2019 - 11 - 2		2019 - 12 - 12	BJS LLC	200726979	49998

Region	Baseline - Date	Credit - Limit	DUNN LEVEL	Is Open	Invoiced	Customer Payment Terms	Total Open Amount - USD	USD - CURRENCY	Total Open Amount	Invoice Currency	Pay term	Due - Date	Posting - Date	Business Year	Payment - Date	Customer Name	Customer Number	Business Code
SOUTHEAST	2019-01-02	20000	1	0	1928576e+09	NAM4	1790300	USD	1790	USD	19	2019-01-15	2019-01-05	2019	2019-01-15	DEC corp	200020431	U0001
49999																		


```
DateParseError Traceback (most recent call last) Cell In[57], line 1 ----> 1
df_receivables_data['Payment_Date'] =
pandas.to_datetime(df_receivables_data['Payment_Date']) 2 df_receivables_data.info()
```

```
File ~/Documents/КПК/Чемпионат/2026/машинное
обучение/Юниоры/work/env/lib/python3.12/site-packages/pandas/core/tools/datetimes.py:1
036, in to_datetime(arg, errors, dayfirst, yearfirst, utc, format, exact, unit, origin, cache) 1034
result = arg.tz_localize("utc") 1035 elif isinstance(arg, ABCSeries): -> 1036 cache_array =
_maybe_cache(arg, format, cache, convert_listlike) 1037 if not cache_array.empty: 1038
result = arg.map(cache_array)
```

```
File ~/Documents/КПК/Чемпионат/2026/машинное
обучение/Юниоры/work/env/lib/python3.12/site-packages/pandas/core/tools/datetimes.py:2
54, in _maybe_cache(arg, format, cache, convert_listlike) 252 unique_dates = unique(arg)
253 if len(unique_dates) < len(arg): --> 254 cache_dates = convert_listlike(unique_dates,
format) 255 # GH#45319 256 try:
```

```
File ~/Documents/КПК/Чемпионат/2026/машинное
обучение/Юниоры/work/env/lib/python3.12/site-packages/pandas/core/tools/datetimes.py:4
37, in _convert_listlike_datetimes(arg, format, name, utc, unit, errors, dayfirst, yearfirst,
exact) 434 if format is not None and format != "mixed": 435 return
_array_strptime_with_fallback(arg, name, utc, format, exact, errors) --> 437 result, tz_parsed
= objects_to_datetime64( 438 arg, ... File pandas/_libs/tslibs/parsing.pyx:307, in
pandas._libs.tslibs.parsing.parse_datetime_string()
```

```
File pandas/_libs/tslibs/parsing.pyx:218, in
pandas._libs.tslibs.parsing._parse_delimited_date()
```

DateParseError: Invalid date specified (1/0) Output is truncated. View as a scrollable element or open in a text editor. Adjust cell output settings...

Данная ошибка говорит о том, что в форматированном столбце есть несоответствие формату даты. Нужно исправить.

In [167]:

```
df_receivables_data['Payment_Date'] =
```



```
pandas.to_datetime(df_receivables_data['Payment_Date'], errors='coerce',
format="%y%m%d")
```

Проверяем

In [168]:

```
df_receivables_data.info()
```

```
<class 'pandas.DataFrame'>
Index: 49928 entries, 0 to 49999
Data columns (total 19 columns):
#   Column                Non-Null Count  Dtype
---  -
0   Business Code         49928 non-null  str
1   Customer Number       49928 non-null  object
2   Customer Name         49928 non-null  str
3   Payment_Date          39942 non-null  datetime64[us]
4   Business Year         49928 non-null  int64
5   Posting_Date          49928 non-null  datetime64[us]
6   Due_Date              49928 non-null  datetime64[us]
7   Payterm               49928 non-null  int64
8   Invoice Currency       49928 non-null  str
9   Total Open Amount     49928 non-null  int64
10  USD_CURRENNCY         49928 non-null  str
11  Total Open Amount_USD 49928 non-null  float64
12  Customer Payment Terms 49928 non-null  str
13  Invoice ID             49922 non-null  float64
14  Is Open               49928 non-null  int64
15  DUNNLEVEL             49928 non-null  int64
16  Credit_limit          49928 non-null  int64
17  Baseline_Date         49928 non-null  datetime64[us]
18  Region                49928 non-null  str
dtypes: datetime64[us](4), float64(2), int64(6), object(1), str(6)
```

memory usage: 7.6+ MB

In [169]:

df_receivables_data

Out[169]:

Region	BaselLine_Date	CreditLimit	DUNNLEVEL	IsOpen	InvoiceID	CustomerPaymentTerms	TotalOpenAmount_USD	USD_CURRENCY	TotalOpenAmount	InvoiceCurrency	Payterm	DueDate	PostingDate	BusinessYear	PaymentDate	CustomerName	CustomerNumber	BusinessCode
WEST	2020-01-2	500000	0	0	1930438e	NAH4	54273.28	USD	54273	USD	15	2020-02-1	2020-01-2		2020	WAL-MARCO	20076962	U001

	B u s i n e s s C o d e	C u s t o m e r N u m b e r	C u s t o m e r N a m e	P a y m e n t _ D a t e	B u s i n e s s Y e a r	P o s t i n g _ D a t e	D u e _ D a t e	P a y t e r m	I n v o i c e C u r r e n c y	T o t a l O p e n A m o u n t	U S D - C U R R E N C Y	T o t a l O p e n A m o u n t - U S D	C u s t o m e r P a y m e n t T e r m s	I n v o i c e I D	I s O p e n	D U N N L E V E L	C r e d i t _ L i m i t	B a s e l i n e _ D a t e	R e g i o n
1	U0001	200980828	BENE	2019-08-08	2019	2019-07-22	2019-08-11	20	USD	79657	USD	79656.600	NAD1	.929646e+09	0	2	50000	2019-07-20	MIDWEST
2	U0001	2000792	MDV / tr	2019-11-	2019	2019-11-00	2019-11-00	15	USD	2254	USD	2253.8	NAA8	.9298	0	3	10000	2019-0	NORTHE

	B u s i n e s s C o d e	C u s t o m e r N u m b e r	C u s t o m e r N a m e	P a y m e n t _ D a t e	B u s i n e s s Y e a r	P o s t i n g _ D a t e	D u e _ D a t e	P a y t e r m	I n v o i c e C u r r e n c y	T o t a l O p e n A m o u n t	U S D - C U R R E N C Y	T o t a l O p e n A m o u n t _ U S D	C u s t o m e r P a y m e n t T e r m s	I n v o i c e I D	I s O p e n	D U N N L E V E L	C r e d i t _ L i m i t	B a s e l i n e _ D a t e	R e g i o n
3	C A 0 2	7 3 4	S Y S C I l c	2 - 3 0	2 0 2 0	9 - 1 4	9 - 2 9	1 1	C A D	3 3 0 0	U S D	2 4 4 1 . 7 7 8	C A 1 0	7 4 e + 0 9	1	2	1 0 0 0 0 0	9 - 1 4	A S T
4	U 0 0	2 0 0	W A L	2 0 1	2 0 1	2 0 1	2 0 1	1 5	U S	3 3 1	U S	3 3 1	N A H	1 . 9	0	3	1 0 0	2 0 1	W E S

B u s i n e s s C o d e	C u s t o m e r N u m b e r	C u s t o m e r N a m e	P a y m e n t _ D a t e	B u s i n e s s Y e a r	P o s t i n g _ D a t e	D u e _ D a t e	P a y t e r m	I n v o i c e C u r r e n c y	T o t a l O p e n A m o u n t	U S D - C U R R E N C Y	T o t a l O p e n A m o u n t _ U S D	C u s t o m e r P a y m e n t T e r m s	I n v o i c e I D	I s O p e n	D U N L E V E L	C r e d i t _ L i m i t	B a s e L i n e _ D a t e	R e g i o n
1	769623	- M A R f o u n d a t i o n	9-11-25	9	9-11-13	9-11-28		D	33	D	33.290	4	30148e+09			000	9-11-10	T
.
4999	U001	20056	C O c o r	N a T	2020-	2020-	15	U S D	3188	U S D	3187.	N A A 8	1 . 9 3 0	1	3	10000	2020-	N O R T H

Region	Baseline - Date	Credit - Limit	DUNN LEVEL	Is Open	Invoice ID	Customer Payment Terms	Total Open Amount - USD	USD - CURRENCY	Total Open Amount	Invoice Currency	Pay term	Due - Date	Posting - Date	Business Year	Payment - Date	Customer Name	Customer Number	Business Code
EAST	04 - 21	0			797e + 09		860					05 - 06	04 - 21			poration	1861	5
WEST	2019 - 08 - 11	100000	2	0	1929744e + 09	NAH4	6766.540	USD	6767	USD	15	2011 - 08 - 30	2011 - 08 - 30	2019	2019 - 03	WAL - MARCHCO	200769623	49996
SO	20	10	1	0	1.	NA	61	US	61	US	1	20	20	20	20	SA	20	49

Region	Baseline - Date	Credit - Limit	DUNN LEVEL	Is Open	Invoice ID	Customer Payment Terms	Total Open Amount - USD	USD - CURRENCY	Total Open Amount	Invoice Currency	Payer term	Due - Date	Posting - Date	Business Year	Payment - Date	Customer Name	Customer Number	Business Code
UTHEAST	20 - 02 - 16	0000			930537e+09	A8	20 . 860	D	21	D	5	20 - 03 - 05	20 - 02 - 19	20	20 - 03 - 05	FEW Associates	0772595	997
MIDWEST	2019 - 11 - 2	100000	0	0	1 . 930199e+	NAA8	63 . 480	USD	63	USD	15	2019 - 12 - 1	2019 - 11 - 2	20	2019 - 12 - 1	BJS LLC	00726979	4998

Region	Baseline Date	Credit Limit	DUNLEVEL	Is Open	Invoice ID	Customer Payment Terms	Total Open Amount - USD	USD - CURRENCY	Total Open Amount	Invoice Currency	Pay term	Due Date	Posting Date	Business Year	Payment Date	Customer Name	Customer Number	Business Code
SOUTHEAST	2019-01-02	20000	1	0	1928576e+09	NAM4	1790.300	USD	1790	USD	19	2019-01-24	2019-01-05	2019	2019-01-15	DEC corp	200020431	U001
49999																		

49928 rows × 19 columns

Исследуем выбросы

In [170]:

```
df_receivables_data.describe()
```

Out[170]:

	Payment_Date	Business Year	Posting_Date	Due_Date	Payterm	Total Open Amount	Total Open Amount_USD	Invoice ID	Is Open	DUNN LEVEL	Credit_limit	Baseline_Date
count	39942	49928.0000	49928	49928	49928.0000	49928.0000	49928.0000	4.99200e+04	49928.0000	49928.0000	49928.0000	49928
mean	2019-08-10 08:25:46.642631	2019-09-11 20:44:15.151418	2019-09-11 20:44:15.151418	2019-09-28 23:27:437910	17.112642	32331.848161	31186.62573	2.011251e+09	0.20008	1.503725	78110.679378	2019-09-20:54:45.050472
min	2019-	2019.	2018-	2018-	-50.00	1.00	0.720	1.928	0.00	0.00	0.00	2018-

	Pa ym ent _D ate	Bu sin es _Y e ar	Pos ti ng _D ate	Du e _D ate	Pa ym ent	Tot al Op en Am ou nt	Tot al Op en Am ou nt _US D	Inv oi ce ID	Is Op en	DU NN LE VEL	Cr edi t _li mit	Bas eline _D ate
n	01-03-00:00:00	00-00-00	12-30-00:00:00	12-24-00:00:00	0000	000	000	502e+09	000	000	000	12-26-00:00:00
25%	2019-05-01-00:00:00	2019-05-06-00:00:00	2019-05-06-00:00:00	2019-05-24-00:00:00	15.0000	4924.0000	4776.51070	1.929342e+09	0.0000	1.0000	5000.0000	2019-05-04-00:00:00
50%	2019-08-07-00:00:00	2019-09-09-00:00:00	2019-09-26-00:00:00	2019-09-26-00:00:00	15.0000	17604.0000	17296.3150	1.929964e+09	0.0000	2.0000	1000.0000	2019-09-07-00:00:00
75%	2019-11-15-00:00:00	2020-01-31-00:00:00	2020-02-16-00:00:00	2020-02-16-00:00:00	15.0000	47124.0000	46131.86750	1.930619e+09	0.0000	3.0000	1000.0000	2020-01-29-00:00:00
max	2020-01-01-00:00:00	2020-01-01-00:00:00	2020-01-01-00:00:00	2020-01-01-00:00:00	120.0000	6685.0000	6685.0000	2.960.0000	1.0000	3.0000	5000.0000	2020-01-01-00:00:00

	Pa ym ent _D ate	Bu sin es _Y e ar	Pos ti ng _D ate	Du e_ Dat e	Pa yte rm	Tot al Op en Am ou nt	Tot al Op en Am ou nt_ US D	Inv oic e ID	Is Op en	DU NN LE VE L	Cr edi t_ li mit	Bas eline _D ate
	05- 22 00: 00: 00	00 00 00	05- 22 00: 00: 00	07- 10 00: 00: 00	00 00 0	93. 00 00 00	93. 36 00 0	63 6e +0 9	00 0	00 0	00. 00 00 00	05- 19 00: 00: 00
std	NaN	0.4 60 69 1	NaN	NaN	10. 26 97 22	39 21 2.5 65 29 6	36 61 0.7 71 50	2.7 64 91 8e +0 8	0.4 00 01 0	1.1 18 05 7	42 98 8.1 33 92 1	NaN

Вывод: Столбец 'Payterm' (Срок выплат) содержит минимальное отрицательное значение, что не может быть. Высним сколько всего отрицательных значений в этом столбце

In [171]:

```
df_receivables_data[df_receivables_data['Payterm'] < 0]
```

Out[171]:

[illegible]

[illegible]

[illegible]

Region	Baseline - Date	Credit - Limit	DUNN LEVEL	Is Open	Invoice ID	Customer Payment Terms	Total Open Amount - USD	USD - CURRENCY	Total Open Amount	Invoice Currency	Pay term	Due - Date	Posting - Date	Business Year	Payment - Date	Customer Name	Customer Number	Business Code
SOUTHEAST	2019-05-24	1000000	0	0	1929376e+09	NAM1	12776.45	USD	12776	USD	- 2	2019-05-23	2019-05-25	2019	2019-05-25	DECco	200020431	48093
SOUTHWES	2019-01-	200000	0	0	1928696	NAM2	9773.84	USD	9774	USD	- 2	2019-01-	2019-01-	2019	2019-01-	DECcorp	200041683	48360

Region	Baseline - Date	Credit - Limit	DUNN LEVEL	Is Open	Invoice ID	Customer Payment Terms	Total Open Amount - USD	USD - CURRENCY	Total Open Amount	Invoice Currency	Pay term	Due - Date	Posting - Date	Business Year	Payment - Date	Customer Name	Customer Number	Business Code
T	26				e + 09							26	28	28	28		7	
MIDWEST	2020 - 01 - 26	200000	0	0	1930444e + 09	NAM2	319060	USD	3191	USD	- 1	2020 - 01 - 26	2020 - 01 - 27	2020 - 01 - 29		DEC trust	200803720	49400
MIDWE	2020 -	50000	3	0	1930	NAM2	8518	USD	8518	USD	- 1	2020 -	2020 -	2020 -		DEC sy	20080	4958

Region	Baseline - Date	Credit - Limit	DUNNLEVEL	Is Open	Invoice ID	Customer Payment Terms	Total Open Amount - USD	USD - CURRENCY	Total Open Amount	Invoice Currency	Pay term	Due - Date	Posting - Date	Business Year	Payment - Date	Customer Name	Customer Number	Business Code
ST	01-27	0			445e+09		44					01-26	01-27	0	01-27	stems	3720	9
SOUTHWEST	2019-08-25	20000	1	0	1929802e+09	NAM1	18724.67	USD	18725	USD	-5	2019-08-28	2019-08-28	2	2019-08-28	DEC Aus	200592182	49592

Всего таких аномальных строк 140. Удалим их

In [172]:

```
df_receivables_data = df_receivables_data[df_receivables_data["Payterm"] > 0]
```

In [173]:

```
df_receivables_data
```

Out[173]:

		Region		Baseline - Date	Credit - Limit	DUNN LEVEL	Is Open	Invoice ID	Customer Payment Terms	Total Open Amount - USD	USD - CURRENCY	Total Open Amount	Invoice Currency	Payer term	Due - Date	Posting - Date	Business Year	Payment - Date	Customer Name	Customer Number	Business Code
0	U0001	WEST	1	2020-01-23	50000	0	0	1930438e+09	NAH4	54273.280	USD	54273	USD	15	2020-02-10	2020-01-26	2020-02-01	2020-02-01	WAL-MAR Corp	200769623	U0001
				2019-07-28	50000	2	0	1929656.600	NAD1	79656.600	USD	79657	USD	20	2019-07-22	2019-07-22	2019-08-08	2019-08-08	BENE	200980828	U0001
1	U0001	MIDWEST	1	2019-07-20	50000	2	0	1929656.600	NAD1	79656.600	USD	79657	USD	20	2019-07-22	2019-07-22	2019-08-08	2019-08-08	BENE	200980828	U0001
				2019-07-20	50000	2	0	1929656.600	NAD1	79656.600	USD	79657	USD	20	2019-07-22	2019-07-22	2019-08-08	2019-08-08	BENE	200980828	U0001

[illegible]

[illegible]

Region	Baseline - Date	Credit - Limit	DUNN LEVEL	Is Open	Invoice ID	Customer Payment Terms	Total Open Amount - USD	USD - CURRENCY	Total Open Amount	Invoice Currency	Pay term	Due - Date	Posting - Date	Business Year	Payment - Date	Customer Name	Customer Number	Business Code
.
NORTHEAST	2020-04-21	100000	3	1	1930797e+09	NA A8	3187.860	USD	3188	USD	15	2020-05-06	2020-04-21	2020	NAT	CO Corporation	200561861	49995
WEST	2019-08	100000	2	0	192974	NA H4	6766.54	USD	6767	USD	15	2019-08	2019-08	2019	WAL - MAR	WAL - MAR	2007696	49996

49997

BUSINESS CODE		CUSTOMER NUMBER		CUSTOMER NAME		PAYMENT DATE		BUSINESS YEAR		POSTING DATE		DUE DATE		PAYER		INVOICE CURRENCY		TOTAL OPEN AMOUNT - USD		USD - CURRENCY		CUSTOMER PAYMENT TERMS		INVOICED		IS OPEN		DUNN LEVEL		CREDIT LIMIT		BASELINE DATE		REGION	
200772595		200772595		SAFEWASSOCIATES		2020-03-05		2020-02-01		2020-03-05		2020-03-05		15		USD		6120.860		USD		NA A 8		1.930537e+09		0		1		1000000		2020-02-16		SOUTHEAST	
200772595		200772595		SAFEWASSOCIATES		2020-03-05		2020-02-01		2020-03-05		2020-03-05		15		USD		6120.860		USD		NA A 8		1.930537e+09		0		1		1000000		2020-02-16		SOUTHEAST	
200772595		200772595		SAFEWASSOCIATES		2020-03-05		2020-02-01		2020-03-05		2020-03-05		15		USD		6120.860		USD		NA A 8		1.930537e+09		0		1		1000000		2020-02-16		SOUTHEAST	
200772595		200772595		SAFEWASSOCIATES		2020-03-05		2020-02-01		2020-03-05		2020-03-05		15		USD		6120.860		USD		NA A 8		1.930537e+09		0		1		1000000		2020-02-16		SOUTHEAST	

		Region		Baseline - Date	Credit - Limit	DUNN LEVEL	Is Open	Invoice ID	Customer Payment Terms	Total Open Amount - USD	USD - CURRENCY	Total Open Amount	Invoice Currency	Payer term	Due - Date	Posting - Date	Business Year	Payment - Date	Customer Name	Customer Number	Business Code
		M I D W E S T		2019-11-27	1000000	0	0	19301999e+09	NA A 8	63.480	USD	63	USD	15	2019-11-22	2019-11-27	2019	2019-11-22	B J ' S I l l c	200726979	49998
		S O U T H E A S T		2019-01-02	200000	1	0	1928576e+0	N A M 4	1790.300	USD	1790	USD	19	2019-01-24	2019-01-05	2019	2019-01-15	D E C c o r p	200020431	49999

Region	Baseline_Date	Credit_Limit	DUNNLEVEL	IsOpen	InvoiceID	CustomerPaymentTerms	TotalOpenAmount_USD	USD_CURRENCY	TotalOpenAmount	InvoiceCurrency	PayerTerm	Due_Date	Posting_Date	BusinessYear	Payment_Date	CustomerName	CustomerNumber	BusinessCode
--------	---------------	--------------	-----------	--------	-----------	----------------------	---------------------	--------------	-----------------	-----------------	-----------	----------	--------------	--------------	--------------	--------------	----------------	--------------

9

49077 rows × 19 columns

Вывод: Как видим объем датафрейма уменьшился с 49928 до 49077

Анализируем файл suppliers_data.xls

In [174]:

```
df_suppliers_data = pandas.read_excel("a/xls/suppliers_data.xls")
```

```
df_suppliers_data.info()
```

```
<class 'pandas.DataFrame'>
RangeIndex: 500 entries, 0 to 499
Data columns (total 5 columns):
#   Column          Non-Null Count  Dtype
---  ---
0   Supplier ID      500 non-null   str
1   Supplier Name     500 non-null   str
2   Payment Terms    500 non-null   str
3   Vendor Type       500 non-null   str
4   Supplier Category 500 non-null   str
dtypes: str(5)
memory usage: 19.7 KB
```

In [175]:

```
df_suppliers_data
```

Out[175]:

	Supplier ID	Supplier Name	Payment Terms	Vendor Type	Supplier Category
0	S-281	Roth-Sanchez	Net 60	Domestic	Raw Material
1	S-438	Peterson Inc	Net 60	Domestic	Raw Material

	Supplier ID	Supplier Name	Payment Terms	Vendor Type	Supplier Category
2	S-480	Morton, Newman and Baker	Net 90	Domestic	Services
3	S-148	Evans Inc	Net 30	Domestic	Utility
4	S-8	Hart Ltd	Net 90	International	Taxes
...
495	S-345	Terrell-Wyatt	Net 60	Domestic	Raw Material
496	S-132	Porter-Anderson	Net 60	Domestic	Raw Material
497	S-451	Brown-Fisher	Net 30	Domestic	Raw Material
498	S-235	Warren and Sons	Net 90	Domestic	Raw Material
499	S-321	Pollard Group	Net 90	Domestic	Raw Material

500 rows × 5 columns

Подсчитываем сумму пустых значений в столбцах датафрейма

In [176]:

```
df_suppliers_data.isna().sum()
```

Out[176]:

```
Supplier ID      0
Supplier Name    0
Payment Terms    0
Vendor Type      0
Supplier Category 0
dtype: int64
```

Считаем дубликаты

In [177]:

```
df_suppliers_data.duplicated().sum()
```

Out[177]:

```
np.int64(0)
```

Выбросы

In [178]:

```
df_customers_data.describe()
```

Out[178]:

Credit Limit

count	4574.000000
mean	89435.942282
std	178507.610874
min	5000.000000
25%	20000.000000
50%	50000.000000
75%	100000.000000
max	1000000.000000

Вывод: Аномалии не обнаружены

Загрузка данных в базу данных

Пример для PostgreSQL

```
connection =  
psycopg2.connect(database="tourism",user="postgres",password="admin",port  
=5432,host="localhost")  
cur = connection.cursor()
```

Импорт необходимых библиотек

In [179]:

```
import sqlite3  
import pandas
```

Используем четыре обработанных датафрейма для загрузки в четыре разных таблицы
(база данных одна)

Создаем таблицу customers на основе датафрейма df_customers_data

In [180]:

```
df_customers_data.info()
```

```
<class 'pandas.DataFrame'>
Index: 4574 entries, 0 to 4976
Data columns (total 5 columns):
#   Column                                Non-Null Count  Dtype
---  ---
0   Customer ID                          4574 non-null   str
1   Customer Name                        4574 non-null   str
2   Customer Payment Terms              4574 non-null   str
3   Address                             4574 non-null   str
4   Credit Limit                         4574 non-null   int64
dtypes: int64(1), str(4)
memory usage: 214.4 KB
```

In [181]:

```
connection = sqlite3.connect("database.db")
cur = connection.cursor()

customers_data = df_customers_data .values.tolist()

sql_create_table_customers = """
CREATE TABLE IF NOT EXISTS customers(
    "Customer ID" TEXT PRIMARY KEY ,
    "Customer_Name" TEXT,
    "Customer_Payment_Terms" TEXT,
    "Address" TEXT,
    "Credit_Limit" INTEGER
);
"""
```

```

cur.execute(sql_create_table_customers)
connection.commit()

for i, row in enumerate(customers_data):
    sql_insert_customers = """
        INSERT INTO customers(
            "Customer ID",
            "Customer_Name",
            "Customer_Payment_Terms",
            "Address",
            "Credit_Limit"
        )
        VALUES (?, ?, ?, ?, ?);
    """
    cur.execute(sql_insert_customers, (row[0], row[1], row[2], row[3], row[4]))
    connection.commit()

connection.close()

```

In [182]:

```
df_payables_data.info()
```

```

<class 'pandas.DataFrame'>
RangeIndex: 20000 entries, 0 to 19999
Data columns (total 15 columns):
#   Column                                     Non-Null Count  Dtype
---  -
0   Invoice Number                             20000 non-null  str
1   Posting Date                               20000 non-null  datetime64[us]
2   Invoice Date                               20000 non-null  datetime64[us]
3   Payment Date                              13988 non-null  datetime64[us]
4   Net Due Date (System Calculated Date)      20000 non-null  datetime64[us]
5   Supplier ID                               20000 non-null  str
6   Invoice Amount                             20000 non-null  int64
7   Fiscal year                               20000 non-null  str
8   Overdue                                    20000 non-null  int64
9   Invoice Status                             20000 non-null  str
10  Spend Category                             20000 non-null  str

```



```
11 Total Outstanding amount          20000 non-null int64
12 Late payment fees                  20000 non-null int64
13 Payterm_n                          20000 non-null int64
14 Vendor_Type                        20000 non-null str
dtypes: datetime64[us](4), int64(5), str(6)
memory usage: 2.3 MB
```

In [183]:

```
connection = sqlite3.connect("database.db")
cur = connection.cursor()

payables_data = df_payables_data .values.tolist()

sql_create_table_payables = """
CREATE TABLE IF NOT EXISTS payables(
    "Invoice Number" TEXT PRIMARY KEY,
    "Posting Date" TIMESTAMP,
    "Invoice Date" TIMESTAMP,
    "Payment Date" TIMESTAMP,
    "Net Due Date (System Calculated Date)" TIMESTAMP,
    "Supplier ID" TEXT,
    "Invoice Amount" INTEGER,
    "Fiscal year" TEXT,
    "Overdue" INTEGER,
    "Invoice Status" TEXT,
    "Spend Category" TEXT,
    "Total Outstanding amount" INTEGER,
    "Late payment fees" INTEGER,
    "Payterm_n" INTEGER,
    "Vendor_Type" TEXT
);
"""

cur.execute(sql_create_table_payables)
connection.commit()

for i, row in enumerate(payables_data):
    sql_insert_payables = """
        INSERT INTO payables(
            "Invoice Number",
            "Posting Date",
            "Invoice Date",
```

```

        "Payment Date",
        "Net Due Date (System Calculated Date)",
        "Supplier ID",
        "Invoice Amount",
        "Fiscal year",
        "Overdue",
        "Invoice Status",
        "Spend Category",
        "Total Outstanding amount",
        "Late payment fees",
        "Payterm_n",
        "Vendor_Type"
    )
    VALUES (?, ?, ?, ?, ?, ?, ?, ?, ?, ?, ?, ?, ?, ?, ?);
"""
    cur.execute(sql_insert_payables, (row[0], str(row[1]), str(row[2]), str(row[3]), str(row[4]),
row[5], row[6], row[7], row[8], row[9], row[10], row[11], row[12], row[13], row[14]))
    connection.commit()

connection.close()

```

In [184]:

```
df_receivables_data.info()
```

```

<class 'pandas.DataFrame'>
Index: 49077 entries, 0 to 49999
Data columns (total 19 columns):
#   Column                Non-Null Count  Dtype
---  ---
0   Business Code          49077 non-null  str
1   Customer Number        49077 non-null  object
2   Customer Name          49077 non-null  str
3   Payment_Date           39195 non-null  datetime64[us]
4   Business Year          49077 non-null  int64
5   Posting_Date           49077 non-null  datetime64[us]
6   Due_Date               49077 non-null  datetime64[us]
7   Payterm                49077 non-null  int64
8   Invoice Currency        49077 non-null  str

```

```

9 Total Open Amount      49077 non-null int64
10 USD_CURRENNCY         49077 non-null str
11 Total Open Amount_USD 49077 non-null float64
12 Customer Payment Terms 49077 non-null str
13 Invoice ID              49077 non-null float64
14 Is Open                 49077 non-null int64
15 DUNNLEVEL              49077 non-null int64
16 Credit_limit            49077 non-null int64
17 Baseline_Date           49077 non-null datetime64[us]
18 Region                  49077 non-null str
dtypes: datetime64[us](4), float64(2), int64(6), object(1), str(6)
memory usage: 7.5+ MB

```

In [185]:

```

connection = sqlite3.connect("database.db")
cur = connection.cursor()

receivables_data = df_receivables_data .values.tolist()

sql_create_table_receivables = """
CREATE TABLE IF NOT EXISTS receivables(
    "Business Code" TEXT,
    "Customer Number" TEXT,
    "Customer Name" TEXT,
    "Payment_Date" TIMESTAMP,
    "Business Year" INTEGER,
    "Posting_Date" TIMESTAMP,
    "Due_Date" TIMESTAMP,
    "Payterm" INTEGER,
    "Invoice Currency" TEXT,
    "Total Open Amount" INTEGER,
    "USD_CURRENNCY" TEXT,
    "Total Open Amount_USD" FLOAT,
    "Customer Payment Terms" TEXT,
    "Invoice ID" FLOAT,
    "Is Open" INTEGER,
    "DUNNLEVEL" INTEGER,
    "Credit_limit" INTEGER,
    "Baseline_Date" TIMESTAMP,
    "Region" TEXT
);
"""

```

```

cur.execute(sql_create_table_receivables)
connection.commit()

for i, row in enumerate(receivables_data):
    sql_insert_receivables = """
        INSERT INTO receivables(
            "Business Code",
            "Customer Number",
            "Customer Name",
            "Payment_Date",
            "Business Year",
            "Posting_Date",
            "Due_Date",
            "Payterm",
            "Invoice Currency",
            "Total Open Amount",
            "USD_CURRENNCY",
            "Total Open Amount_USD",
            "Customer Payment Terms",
            "Invoice ID",
            "Is Open",
            "DUNNLEVEL",
            "Credit_limit",
            "Baseline_Date",
            "Region"
        )
        VALUES (?, ?, ?, ?, ?, ?, ?, ?, ?, ?, ?, ?, ?, ?, ?, ?);
    """

    cur.execute(sql_insert_receivables, (row[0], row[1], row[2], str(row[3]), row[4],
    str(row[5]), str(row[6]), row[7], row[8], row[9], row[10], row[11], row[12], row[13], row[14],
    row[15], row[16], str(row[17]), row[18]))
    connection.commit()

connection.close()

```

In [186]:

```
df_suppliers_data.info()
```

```

<class 'pandas.DataFrame'>
RangeIndex: 500 entries, 0 to 499
Data columns (total 5 columns):
#   Column          Non-Null Count  Dtype
---  ---
0   Supplier ID      500 non-null   str
1   Supplier Name     500 non-null   str
2   Payment Terms     500 non-null   str
3   Vendor Type       500 non-null   str
4   Supplier Category 500 non-null   str
dtypes: str(5)
memory usage: 19.7 KB

```

In [187]:

```

connection = sqlite3.connect("database.db")
cur = connection.cursor()

suppliers_data = df_suppliers_data .values.tolist()

sql_create_table_suppliers = """
CREATE TABLE IF NOT EXISTS suppliers(
    "Supplier ID" TEXT PRIMARY KEY,
    "Supplier Name" TEXT,
    "Payment Terms" TEXT,
    "Vendor Type" TEXT,
    "Supplier Category" TEXT
);
"""

cur.execute(sql_create_table_suppliers)
connection.commit()

for i, row in enumerate(customers_data):
    sql_insert_suppliers = """
        INSERT INTO suppliers(
            "Supplier ID",
            "Supplier Name",
            "Payment Terms",
            "Vendor Type",
            "Supplier Category"
        )
    """

```

```
                VALUES (?, ?, ?, ?, ?);
        """
        cur.execute(sql_insert_suppliers, (row[0], row[1], row[2], row[3], row[4]))
        connection.commit()

connection.close()
```