# Use e-commerce data to analyze and classify customer behavior and implement precision marketing

edit by David Yang 02/14/2023

## import module

In [1]:
```python
import pandas as pd
import matplotlib.pyplot as plt
from datetime import datetime
import seaborn as sns
import warnings
warnings.filterwarnings('ignore')
```

## import data，path is './工作/data.csv'，Encoding format is 'utf-8'

In [4]:
```python
missing_values = ['-','na','none','null','']
test_data = pd.read_csv('E:/风变/数据分析实训营/all_data.csv',na_values = missing_values,
test_data.head(10)
```

Out[4]:

| | 订单号 | 顾客ID | 订单时间 | 付款金额 | |
|---|---|---|---|---|---|
| 0 | e481f51cbdc54678b7cc49136f2d6af7 | 9ef432eb6251297304e76186b10a928d | 2017-10-02 10:56 | 18.12 | 87285b3488457 |
| 1 | e481f51cbdc54678b7cc49136f2d6af7 | 9ef432eb6251297304e76186b10a928d | 2017-10-02 10:56 | 2.00 | 87285b3488457 |
| 2 | e481f51cbdc54678b7cc49136f2d6af7 | 9ef432eb6251297304e76186b10a928d | 2017-10-02 10:56 | 18.59 | 87285b3488457 |
| 3 | 128e10d95713541c87cd1a2e48201934 | a20e8105f23924cd00833fd87daa0831 | 2017-08-15 18:29 | 37.77 | 87285b3488457 |
| 4 | 0e7e841ddf8f8f2de2bad69267ecfbcf | 26c7ac168e1433912a51b924fbd34d34 | 2017-08-02 18:24 | 37.77 | 87285b3488457 |
| 5 | bfc39df4f36c3693ff3b63fcbea9e90a | 53904ddbea91e1e92b2b3f1d09a7af86 | 2017-10-23 23:26 | 44.09 | 87285b3488457 |
| 6 | 6ea2f835b4556291ffdc53fa0b3b95e8 | c7340080e394356141681bd4c9b8fe31 | 2017-11-24 21:27 | 356.12 | be021417a6acb |

| | 订单号 | 顾客ID | 订单时间 | 付款金额 | |
|---|---|---|---|---|---|
| 7 | 82bce245b1c9148f8d19a55b9ff70644 | 388025bec8128ff20ec1a316ed4dcf02 | 2017-04-20 17:15 | 267.80 | a5a0e71a81ae6 |
| 8 | 82bce245b1c9148f8d19a55b9ff70644 | 388025bec8128ff20ec1a316ed4dcf02 | 2017-04-20 17:15 | 267.80 | a5a0e71a81ae6 |
| 9 | 82bce245b1c9148f8d19a55b9ff70644 | 388025bec8128ff20ec1a316ed4dcf02 | 2017-04-20 17:15 | 267.80 | a5a0e71a81ae6 |

In [5]:
```python
test_data.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 115878 entries, 0 to 115877
Data columns (total 6 columns):
 #   Column  Non-Null Count   Dtype
---  ------  --------------   -----
 0   订单号     115878 non-null  object
 1   顾客ID    115878 non-null  object
 2   订单时间    115878 non-null  object
 3   付款金额    115878 non-null  float64
 4   商品ID    115878 non-null  object
 5   商品描述    115878 non-null  object
dtypes: float64(1), object(5)
memory usage: 5.3+ MB
```

# Data Cleansing

In [6]:
```python
missing_value = ['-','na','none','null','inf']
```

In [7]:
```python
test_data.isnull().sum()
```

Out[7]:
```
订单号     0
顾客ID    0
订单时间    0
付款金额    0
商品ID    0
商品描述    0
dtype: int64
```

In [8]:
```python
test_data[test_data.duplicated()]
```

Out[8]:

| | 订单号 | 顾客ID | 订单时间 | 付款金额 | |
|---|---|---|---|---|---|
| 8 | 82bce245b1c9148f8d19a55b9ff70644 | 388025bec8128ff20ec1a316ed4dcf02 | 2017-04-20 17:15 | 267.80 | a5a0e71a |

| | 订单号 | 顾客ID | 订单时间 | 付款金额 | |
|---|---|---|---|---|---|
| 9 | 82bce245b1c9148f8d19a55b9ff70644 | 388025bec8128ff20ec1a316ed4dcf02 | 2017-04-20 17:15 | 267.80 | a5a0e71a |
| 10 | 82bce245b1c9148f8d19a55b9ff70644 | 388025bec8128ff20ec1a316ed4dcf02 | 2017-04-20 17:15 | 267.80 | a5a0e71a |
| 11 | 82bce245b1c9148f8d19a55b9ff70644 | 388025bec8128ff20ec1a316ed4dcf02 | 2017-04-20 17:15 | 267.80 | a5a0e71a |
| 24 | c49be9a11fd13933307cc6a19b03a895 | a972623b3481cbfd95fa776b0067e554 | 2018-05-15 18:54 | 928.68 | 97f1396a |
| ... | ... | ... | ... | ... | ... |
| 115714 | 5020a3db49225f967490d76021c7d13a | 5a8b3e70cb6bfdbc353bcb5ae2b4d4eb | 2018-01-28 23:36 | 188.45 | 3fdb534d |
| 115715 | 5020a3db49225f967490d76021c7d13a | 5a8b3e70cb6bfdbc353bcb5ae2b4d4eb | 2018-01-28 23:36 | 188.45 | 3fdb534d |
| 115716 | 5020a3db49225f967490d76021c7d13a | 5a8b3e70cb6bfdbc353bcb5ae2b4d4eb | 2018-01-28 23:36 | 188.45 | 3fdb534d |
| 115737 | b144e2ac9863ed27bc59dbe4dd2f8773 | 49bc0bacf1f213a2d30e240c648ccb01 | 2017-12-06 14:04 | 99.70 | f83fd2b5 |
| 115781 | 161f105f25baba98c7604aad9b99d9a6 | b9dd6c551bfe1ea46e2ca722708df61d | 2018-03-14 12:26 | 170.60 | 7515ab3 |

11039 rows × 6 columns

In [9]:
```python
test_data1 = test_data.drop_duplicates().reset_index(drop=True)
test_data1[test_data1.duplicated()]
```

Out[9]:

| 订单号 | 顾客ID | 订单时间 | 付款金额 | 商品ID | 商品描述 |
|---|---|---|---|---|---|

In [10]:
```python
test_data1.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 104839 entries, 0 to 104838
Data columns (total 6 columns):
 #   Column   Non-Null Count    Dtype
---  ------   --------------    -----
 0   订单号     104839 non-null   object
 1   顾客ID    104839 non-null   object
```

```
 2    订单时间    104839 non-null  object
 3    付款金额    104839 non-null  float64
 4    商品ID    104839 non-null  object
 5    商品描述    104839 non-null  object
dtypes: float64(1), object(5)
memory usage: 4.8+ MB
```

In [11]:
```python
test_data1.tail()
```

Out[11]:

| | 订单号 | 顾客ID | 订单时间 | 付款金额 | |
|---|---|---|---|---|---|
| 104834 | 0b82d0616f1ad8da15cf967b984b4004 | 986632b40c38f4240caf8608cb01d40d | 2018-08-03 21:35 | 33.69 | 4a24717 |
| 104835 | 2ef4a11b6e24fdfbb43b92cb5f95edff | ee1cfdc92e449920e25d3ca4ab4da4f6 | 2018-07-23 18:35 | 84.63 | 9c313ad |
| 104836 | 2ef4a11b6e24fdfbb43b92cb5f95edff | ee1cfdc92e449920e25d3ca4ab4da4f6 | 2018-07-23 18:35 | 84.63 | eacb1048 |
| 104837 | 2c4ada2e75c2ad41dd93cebb5df5f023 | 363d3a9b2ec5c5426608688ca033292d | 2017-01-26 11:09 | 209.06 | 6c7a0a34 |
| 104838 | bede3503afed051733eeb4a84d1adcc5 | 919570a26efbd068d6a0f66d5c5072a3 | 2017-09-17 16:51 | 115.45 | 8db75af9 |

In [12]:
```python
test_data1['付款金额'].describe()
```

Out[12]:
```
count    104839.000000
mean        158.264636
std         218.993424
min           0.000000
25%          58.370000
50%         102.850000
75%         177.320000
max       13664.080000
Name: 付款金额, dtype: float64
```

### use 3*6 method to remove outlier

plt.hist(test_data1['付款金额'],100,density=True,facecolor='b',alpha=0.8) m = test_data1['付款金额'].mean() std = test_data1['付款金额'].std() plt.axvline(x=m+3*std,color='r') plt.axvline(x=m-3*std,color='r') plt.show()

In [15]:
```python
price_sorted = sorted(test_data1["付款金额"])
threshhold = m+3*std
price_normal = []
price_outlier = []

for price in price_sorted:
```
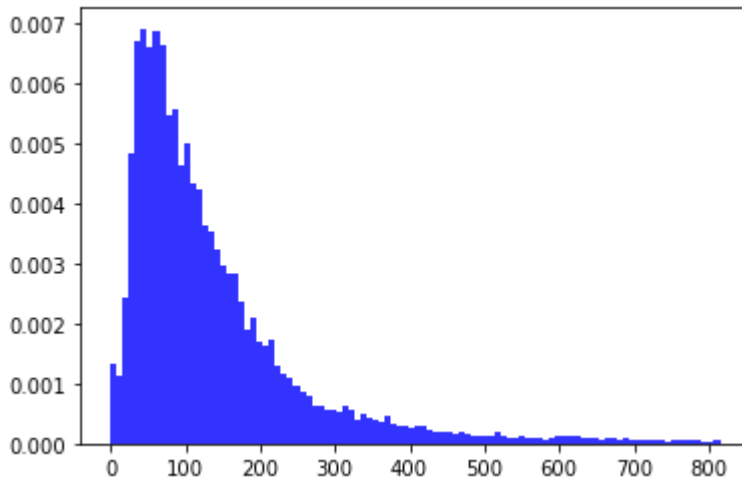
```
    if price<threshhold:
        price_normal.append(price)
    else:
        price_outlier.append(price)

plt.hist(price_normal,100,density=True,facecolor='b',alpha=0.8)
plt.show()
```
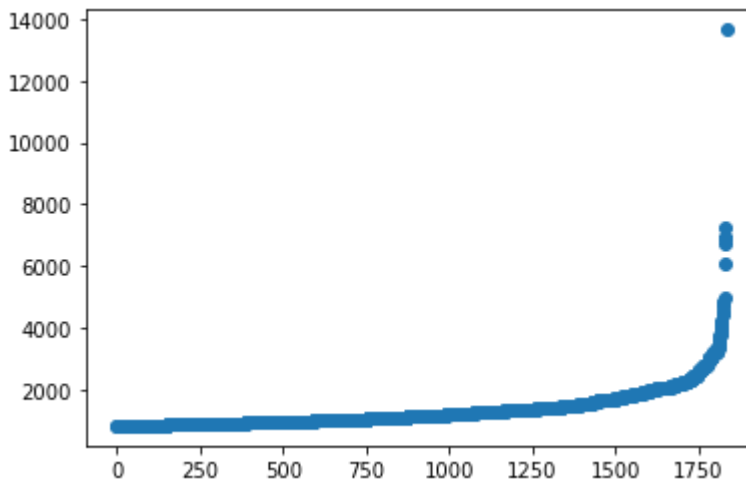


In [16]:
```
len(price_outlier)
```

Out[16]: 1833

In [17]:
```
plt.scatter(range(len(price_outlier)),price_outlier)
plt.show()
```



In [18]:
```
test_data2 = test_data1[test_data1["付款金额"]<m+3*std].reset_index(drop=True)
```
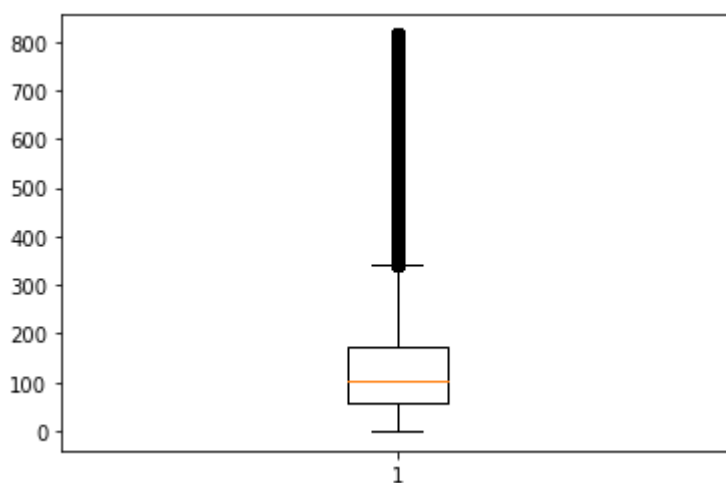
In [19]:
```
test_data2.describe()
```

Out[19]:

| | 付款金额 |
| --- | --- |
| count | 103006.000000 |

| | 付款金额 |
|---|---|
| mean | 137.298398 |
| std | 122.271094 |
| min | 0.000000 |
| 25% | 57.770000 |
| 50% | 100.940000 |
| 75% | 171.780000 |
| max | 814.960000 |

## use 1.5*IQR to analyze based on the 3*6 method to remove outlier

In [21]:
```python
plt.boxplot(test_data2["付款金额"])
plt.show()
```



In [22]:
```python
Q1 = test_data2["付款金额"].quantile(0.25)
Q3 = test_data2["付款金额"].quantile(0.75)
IQR = Q3-Q1
IQR
```
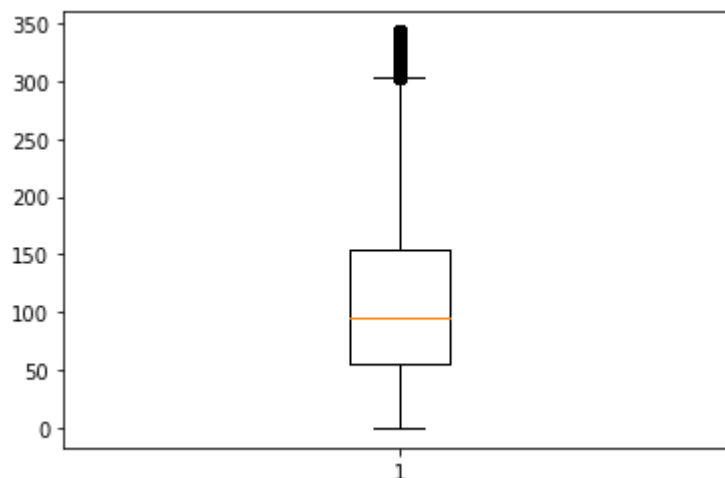
Out[22]:  114.00999999999999

In [23]:
```python
test_data2_normal = test_data2[(test_data2["付款金额"]>Q1-1.5*IQR) & (test_data2["付款金
test_data2_normal.describe()
```

Out[23]:

| | 付款金额 |
|---|---|
| count | 96204.000000 |
| mean | 112.228934 |
| std | 73.282026 |
| min | 0.000000 |
| 25% | 55.240000 |

|  | 付款金额 |
|---|---|
| **50%** | 94.520000 |
| **75%** | 154.200000 |
| **max** | 342.690000 |

In [24]:
```python
plt.boxplot(test_data2_normal["付款金额"])
plt.show()
```



In [25]:
```python
test_data2_normal.info()
```

```
<class 'pandas.core.frame.DataFrame'>
Int64Index: 96204 entries, 0 to 103005
Data columns (total 6 columns):
 #   Column   Non-Null Count  Dtype
---  ------   --------------  -----
 0   订单号      96204 non-null  object
 1   顾客ID     96204 non-null  object
 2   订单时间     96204 non-null  object
 3   付款金额     96204 non-null  float64
 4   商品ID     96204 non-null  object
 5   商品描述     96204 non-null  object
dtypes: float64(1), object(5)
memory usage: 5.1+ MB
```

In [26]:
```python
test_data2_normal.head()
```

Out[26]:

| | 订单号 | 顾客ID | 订单时间 | 付款金额 | |
|---|---|---|---|---|---|
| **0** | e481f51cbdc54678b7cc49136f2d6af7 | 9ef432eb6251297304e76186b10a928d | 2017-10-02 10:56 | 18.12 | 87285b34884572 |
| **1** | e481f51cbdc54678b7cc49136f2d6af7 | 9ef432eb6251297304e76186b10a928d | 2017-10-02 10:56 | 2.00 | 87285b34884572 |
| **2** | e481f51cbdc54678b7cc49136f2d6af7 | 9ef432eb6251297304e76186b10a928d | 2017-10-02 | 18.59 | 87285b34884572 |

| | 订单号 | 顾客ID | 订单时间 | 付款金额 | |
|---|---|---|---|---|---|
| | | | 10:56 | | |
| 3 | 128e10d95713541c87cd1a2e48201934 | a20e8105f23924cd00833fd87daa0831 | 2017-08-15 18:29 | 37.77 | 87285b34884572 |
| 4 | 0e7e841ddf8f8f2de2bad69267ecfbcf | 26c7ac168e1433912a51b924fbd34d34 | 2017-08-02 18:24 | 37.77 | 87285b34884572 |

◄ ▬▬▬▬▬▬▬▬▬▬▬▬▬▬▬▬▬▬▬▬▬ ►

# data wrangling

In [27]:
```python
test_data2_normal['订单时间'] = test_data2_normal['订单时间'].astype('datetime64')
test_data2_normal.info()
```

```
<class 'pandas.core.frame.DataFrame'>
Int64Index: 96204 entries, 0 to 103005
Data columns (total 6 columns):
 #   Column  Non-Null Count  Dtype
---  ------  --------------  -----
 0   订单号     96204 non-null  object
 1   顾客ID    96204 non-null  object
 2   订单时间    96204 non-null  datetime64[ns]
 3   付款金额    96204 non-null  float64
 4   商品ID    96204 non-null  object
 5   商品描述    96204 non-null  object
dtypes: datetime64[ns](1), float64(1), object(4)
memory usage: 5.1+ MB
```

In [28]:
```python
test_data2_normal['year'] = test_data2_normal['订单时间'].dt.year
test_data2_normal['month'] = test_data2_normal['订单时间'].dt.month
test_data2_normal['day'] = test_data2_normal['订单时间'].dt.day
```

In [29]:
```python
test_data2_normal.head()
```

Out[29]:

| | 订单号 | 顾客ID | 订单时间 | 付款金额 | |
|---|---|---|---|---|---|
| 0 | e481f51cbdc54678b7cc49136f2d6af7 | 9ef432eb6251297304e76186b10a928d | 2017-10-02 10:56:00 | 18.12 | 87285b348845 |
| 1 | e481f51cbdc54678b7cc49136f2d6af7 | 9ef432eb6251297304e76186b10a928d | 2017-10-02 10:56:00 | 2.00 | 87285b348845 |
| 2 | e481f51cbdc54678b7cc49136f2d6af7 | 9ef432eb6251297304e76186b10a928d | 2017-10-02 10:56:00 | 18.59 | 87285b348845 |
| 3 | 128e10d95713541c87cd1a2e48201934 | a20e8105f23924cd00833fd87daa0831 | 2017-08-15 | 37.77 | 87285b348845 |

| | 订单号 | 顾客ID | 订单时间 | 付款金额 | |
|---|---|---|---|---|---|
| | | | 18:29:00 | | |
| 4 | 0e7e841ddf8f8f2de2bad69267ecfbcf | 26c7ac168e1433912a51b924fbd34d34 | 2017-08-02 18:24:00 | 37.77 | 87285b348845 |

```
In [31]: test_data2_normal.to_csv('E:/风变/数据分析实训营/cleansing_data.csv',encoding = 'utf-8-si
```

# data analyze

```
In [2]: df = pd.read_csv('E:/风变/数据分析实训营/cleansing_data.csv')
        df.columns = ['order_id','cust_id','order_time','order_payment','pro_id','pro_describe'
        df.head()
```

Out[2]:

| | order_id | cust_id | order_time | order_payment | |
|---|---|---|---|---|---|
| 0 | e481f51cbdc54678b7cc49136f2d6af7 | 9ef432eb6251297304e76186b10a928d | 2017-10-02 10:56:00 | 18.12 | 8 |
| 1 | e481f51cbdc54678b7cc49136f2d6af7 | 9ef432eb6251297304e76186b10a928d | 2017-10-02 10:56:00 | 2.00 | 8 |
| 2 | e481f51cbdc54678b7cc49136f2d6af7 | 9ef432eb6251297304e76186b10a928d | 2017-10-02 10:56:00 | 18.59 | 8 |
| 3 | 128e10d95713541c87cd1a2e48201934 | a20e8105f23924cd00833fd87daa0831 | 2017-08-15 18:29:00 | 37.77 | 8 |
| 4 | 0e7e841ddf8f8f2de2bad69267ecfbcf | 26c7ac168e1433912a51b924fbd34d34 | 2017-08-02 18:24:00 | 37.77 | 8 |

## increase weekday

```
In [3]: df['weekday'] = pd.to_datetime(df['order_time']).dt.weekday
        df.head()
```

Out[3]:

| | order_id | cust_id | order_time | order_payment | |
|---|---|---|---|---|---|
| 0 | e481f51cbdc54678b7cc49136f2d6af7 | 9ef432eb6251297304e76186b10a928d | 2017-10-02 10:56:00 | 18.12 | 8 |
| 1 | e481f51cbdc54678b7cc49136f2d6af7 | 9ef432eb6251297304e76186b10a928d | 2017-10-02 10:56:00 | 2.00 | 8 |
| 2 | e481f51cbdc54678b7cc49136f2d6af7 | 9ef432eb6251297304e76186b10a928d | 2017-10-02 10:56:00 | 18.59 | 8 |
| 3 | 128e10d95713541c87cd1a2e48201934 | a20e8105f23924cd00833fd87daa0831 | 2017-08-15 18:29:00 | 37.77 | 8 |

| | order_id | cust_id | order_time | order_payment |
|---|---|---|---|---|
| **4** | 0e7e841ddf8f8f2de2bad69267ecfbcf | 26c7ac168e1433912a51b924fbd34d34 | 2017-08-02 18:24:00 | 37.77 |

### remove the data in 2016

In [4]:
```python
df.groupby('year')['year'].value_counts()
```

Out[4]:
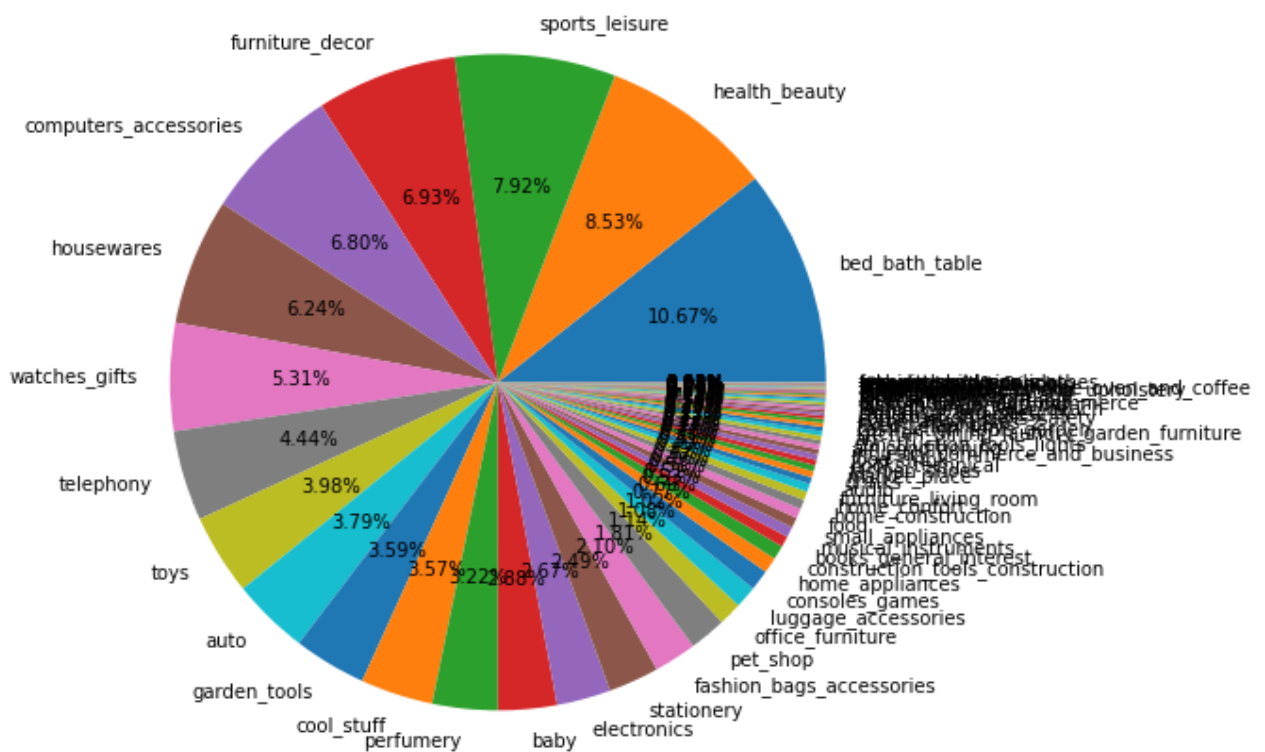```
year  year
2016  2016      304
2017  2017    43687
2018  2018    52213
Name: year, dtype: int64
```

In [5]:
```python
df = df[(df['year']==2017) | (df['year']==2018)]
df.groupby('year')['year'].value_counts()
```

Out[5]:
```
year  year
2017  2017    43687
2018  2018    52213
Name: year, dtype: int64
```

### draw a Pie of product deacribe

In [6]:
```python
ratio_describe = df['pro_describe'].value_counts() / df['pro_describe'].value_counts().
ratio_describe.plot(kind='pie', autopct='%.2f%%',figsize=(8,8),label='')
plt.show()
```

## analyze the sales in different year

In [7]:
```python
plt.rcParams['figure.figsize'] = 10,6
df.groupby(['year','month'])['order_payment'].sum()
```

Out[7]:
```
year  month
2017  1          80203.60
      2         182304.65
      3         280771.10
      4         256491.77
      5         389725.59
      6         339110.83
      7         435607.42
      8         464427.76
      9         459472.04
      10        502989.91
      11        817452.04
      12        620276.42
2018  1         788700.86
      2         716850.69
      3         786641.68
      4         778140.09
      5         764595.04
      6         698249.14
      7         691996.42
      8         706881.42
      9            166.46
Name: order_payment, dtype: float64
```
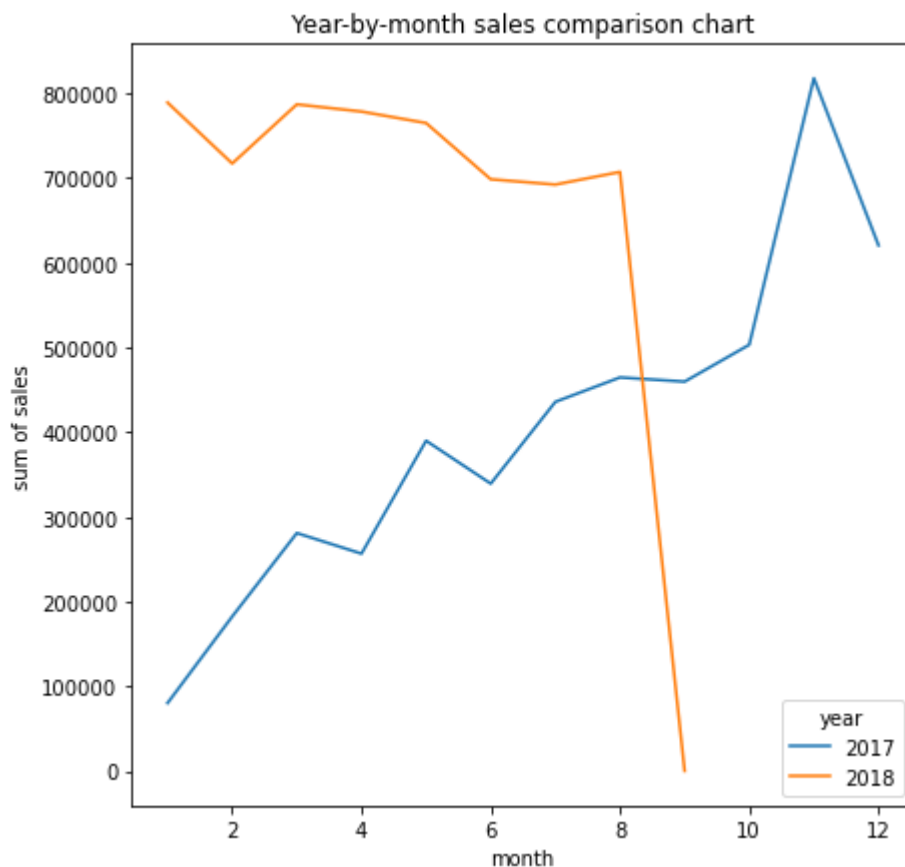
In [8]:
```python
year_month_sales_sum = df.groupby(['year','month'])['order_payment'].sum().unstack(leve
year_month_sales_sum
```

Out[8]:

| year | 2017 | 2018 |
|------|------|------|
| month | | |
| 1 | 80203.60 | 788700.86 |
| 2 | 182304.65 | 716850.69 |
| 3 | 280771.10 | 786641.68 |
| 4 | 256491.77 | 778140.09 |
| 5 | 389725.59 | 764595.04 |
| 6 | 339110.83 | 698249.14 |
| 7 | 435607.42 | 691996.42 |
| 8 | 464427.76 | 706881.42 |
| 9 | 459472.04 | 166.46 |
| 10 | 502989.91 | NaN |
| 11 | 817452.04 | NaN |
| 12 | 620276.42 | NaN |

In [9]:
```python
year_month_sales_sum.plot(kind='line',figsize=(7,7))
plt.xlabel('month')
plt.ylabel('sum of sales')
plt.title('Year-by-month sales comparison chart')
plt.show()
```
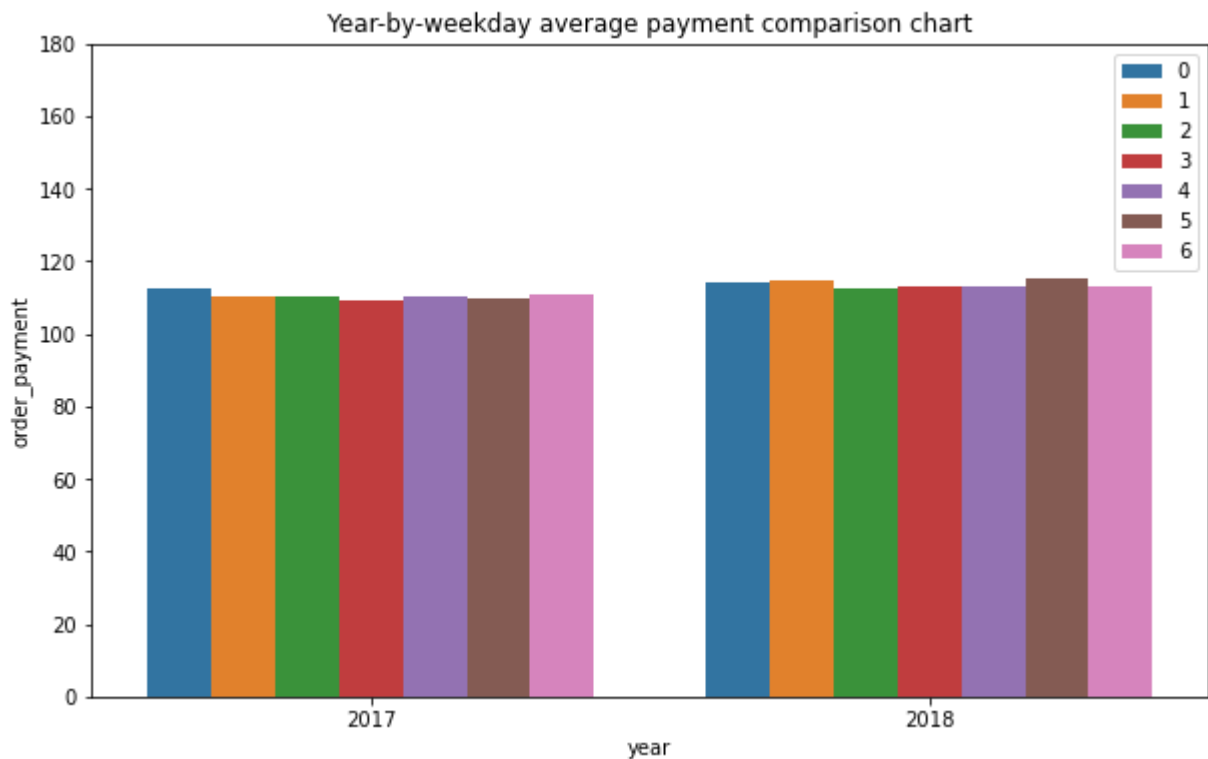


Year-by-month sales comparison chart

use seaborn to draw the means of weekday

In [10]:
```python
import seaborn as sns
```

In [11]:
```python
# to compare the average payment of every weekday in different year
# use hue to class and ci to remove the error interval
sns.barplot(x='year',y='order_payment', data=df, hue='weekday', ci=None)
# extend axis Y
plt.ylim(0,180)
# set the legend
plt.legend(loc='upper right')
plt.title('Year-by-weekday average payment comparison chart')
```

Out[11]:
Text(0.5, 1.0, 'Year-by-weekday average payment comparison chart')

Year-by-weekday average payment comparison chart

## use seaborn to analyze total customers and total payment in every month of 2017

In [12]:
```python
# To compare total of customers in every month, need to combine the same cust_id in the
# after group the data in 2017 by month and cust_id, use agg to get the sum of order_pa
# mean of year just for the chart
df_2017 = df[df['year']== 2017].groupby(['month','cust_id']).agg({'order_payment':'sum'
df_month_customer_2017 = df_2017.reset_index()
df_month_customer_2017
```

Out[12]:

| | month | cust_id | order_payment | year |
|---|---|---|---|---|
| **0** | 1 | 0040b00970e2139e8c43b647c0da5305 | 41.93 | 2017.0 |
| **1** | 1 | 0051337a96842850e1ec728dd158f4b3 | 237.99 | 2017.0 |
| **2** | 1 | 007b7f04a35e02745c23ea706492ca20 | 77.06 | 2017.0 |
| **3** | 1 | 00f3b3a7cd0b6566435090c7fbda03a2 | 57.51 | 2017.0 |
| **4** | 1 | 01a0d45a369a4356ac4652584652109a | 45.86 | 2017.0 |
| **...** | ... | ... | ... | ... |
| **40268** | 12 | ffdb7e488ea7c83b9c1258ee2d3776fa | 85.23 | 2017.0 |
| **40269** | 12 | ffdd933fe636d97903e7a4758faa8c6a | 63.60 | 2017.0 |
| **40270** | 12 | ffe509f377a33554f5a677dcd83e669e | 211.82 | 2017.0 |
| **40271** | 12 | fff675a0d5924b9162b4a1bf410466cd | 75.07 | 2017.0 |
| **40272** | 12 | fff89c8ed4fcf69a823c1d149e429a0b | 44.10 | 2017.0 |

40273 rows × 4 columns

```
In [13]:    # continue to group the data by month and count the total of cumtomer every month
            customer_payment_2017 = df_month_customer_2017.groupby('month').agg({'cust_id':'count',
            customer_payment_2017
```

Out[13]:

| | month | cust_id | order_payment | year |
|---|---|---|---|---|
| 0 | 1 | 684 | 80203.60 | 2017.0 |
| 1 | 2 | 1541 | 182304.65 | 2017.0 |
| 2 | 3 | 2375 | 280771.10 | 2017.0 |
| 3 | 4 | 2123 | 256491.77 | 2017.0 |
| 4 | 5 | 3288 | 389725.59 | 2017.0 |
| 5 | 6 | 2929 | 339110.83 | 2017.0 |
| 6 | 7 | 3642 | 435607.42 | 2017.0 |
| 7 | 8 | 3902 | 464427.76 | 2017.0 |
| 8 | 9 | 3813 | 459472.04 | 2017.0 |
| 9 | 10 | 4126 | 502989.91 | 2017.0 |
| 10 | 11 | 6757 | 817452.04 | 2017.0 |
| 11 | 12 | 5093 | 620276.42 | 2017.0 |

```
In [14]:    customer_payment_2017 = customer_payment_2017.rename(columns={'cust_id':'total_cust','o
            customer_payment_2017
```
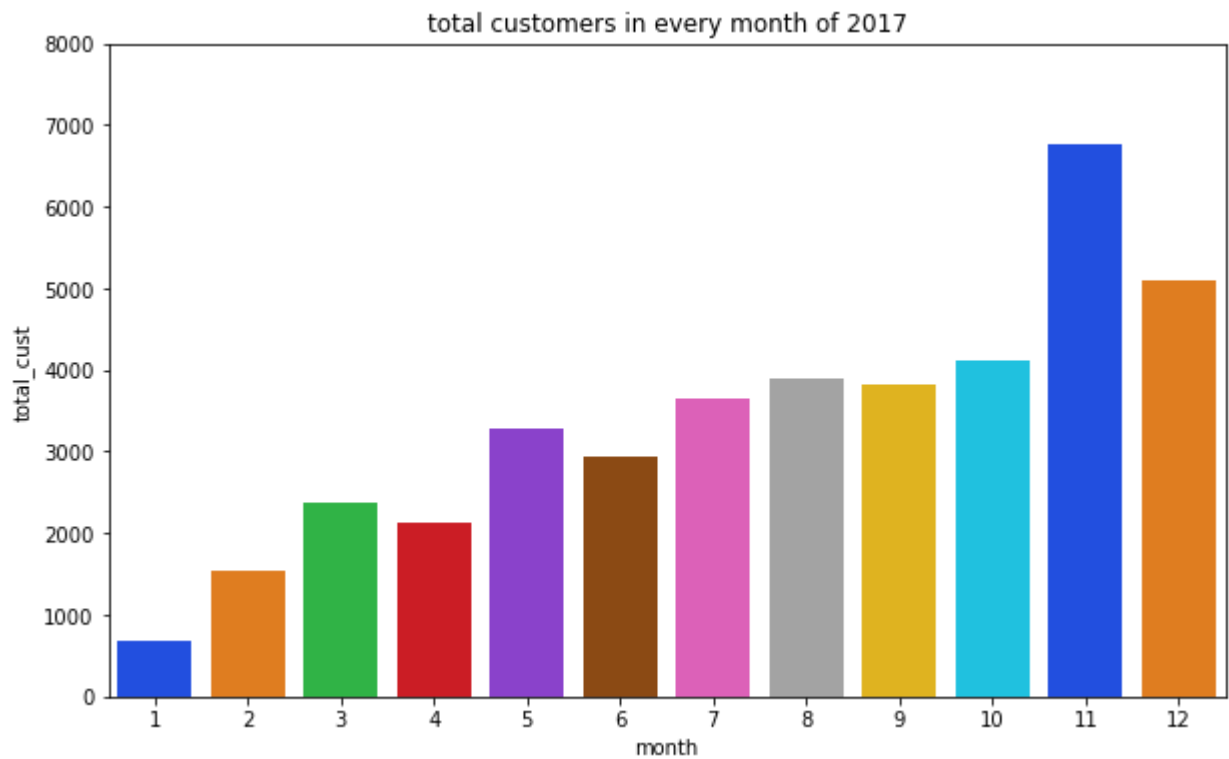
Out[14]:

| | month | total_cust | total_payment | year |
|---|---|---|---|---|
| 0 | 1 | 684 | 80203.60 | 2017.0 |
| 1 | 2 | 1541 | 182304.65 | 2017.0 |
| 2 | 3 | 2375 | 280771.10 | 2017.0 |
| 3 | 4 | 2123 | 256491.77 | 2017.0 |
| 4 | 5 | 3288 | 389725.59 | 2017.0 |
| 5 | 6 | 2929 | 339110.83 | 2017.0 |
| 6 | 7 | 3642 | 435607.42 | 2017.0 |
| 7 | 8 | 3902 | 464427.76 | 2017.0 |
| 8 | 9 | 3813 | 459472.04 | 2017.0 |
| 9 | 10 | 4126 | 502989.91 | 2017.0 |
| 10 | 11 | 6757 | 817452.04 | 2017.0 |
| 11 | 12 | 5093 | 620276.42 | 2017.0 |

```
In [15]:    sns.barplot(x='month',y='total_cust', data=customer_payment_2017, palette='bright')
            # extend axis Y
```
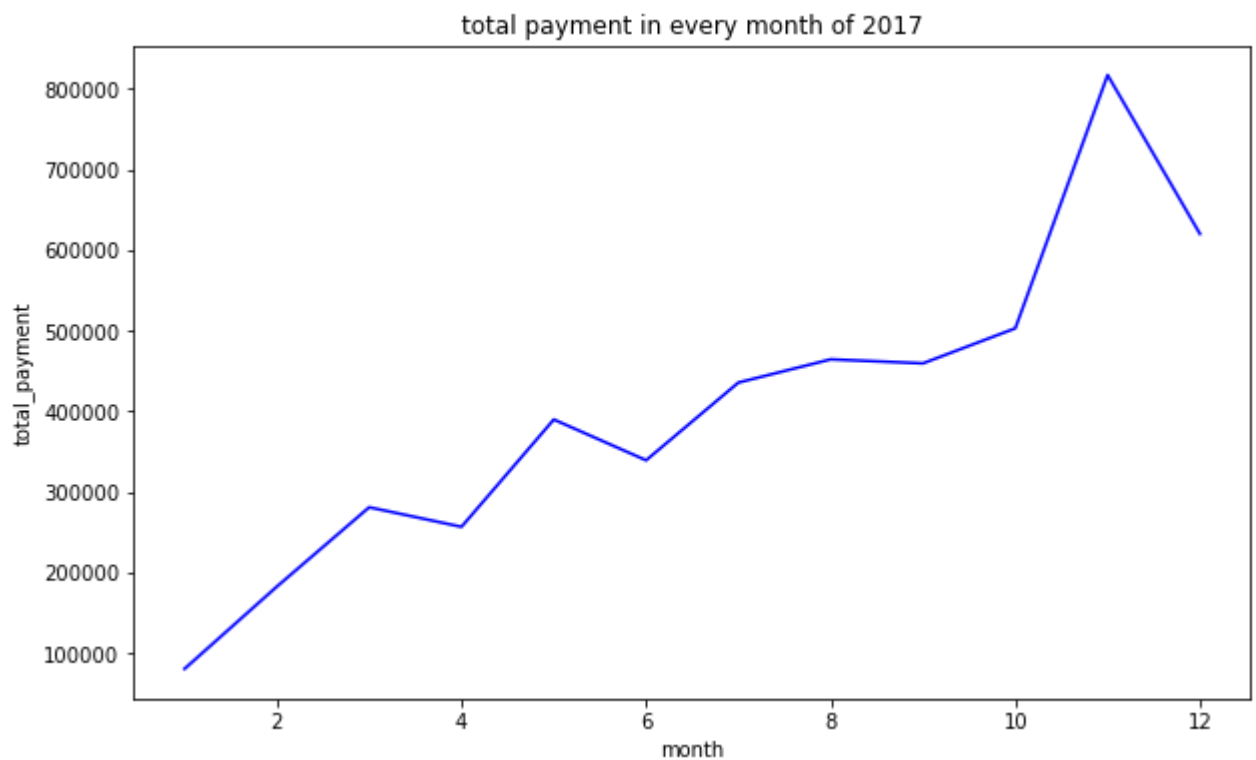
```
plt.ylim(0,8000)
plt.title('total customers in every month of 2017')
```

Text(0.5, 1.0, 'total customers in every month of 2017')

total customers in every month of 2017

```
sns.lineplot(x='month',y='total_payment', data=customer_payment_2017, color='blue')
plt.title('total payment in every month of 2017')
```

Text(0.5, 1.0, 'total payment in every month of 2017')

total payment in every month of 2017

```
In [18]:  df.to_csv('E:/风变/数据分析实训营/analyze_data_1.csv',encoding = 'utf-8-sig',index = Fals
```

```
In [19]:  import pandas as pd
          import matplotlib.pyplot as plt
          from datetime import datetime
          import seaborn as sns
          import warnings
          warnings.filterwarnings('ignore')
```
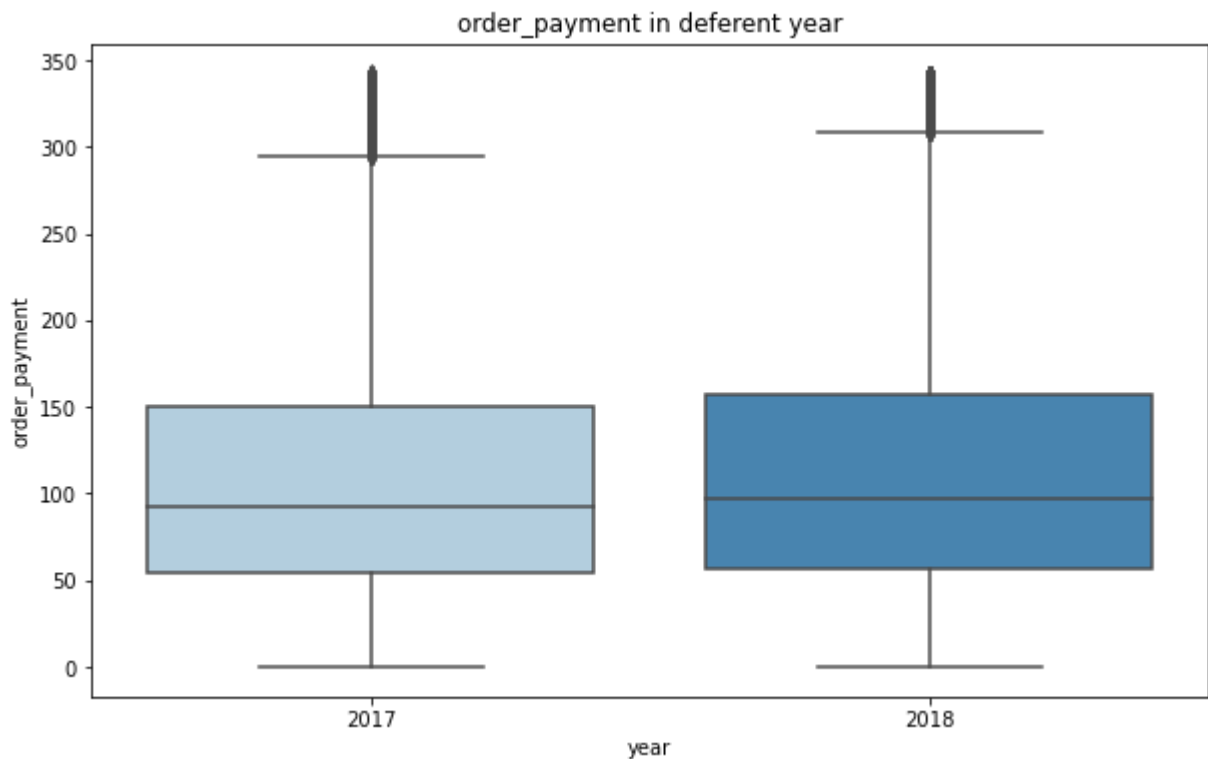
```
In [20]:  df = pd.read_csv('E:/风变/数据分析实训营/analyze_data_1.csv',encoding = 'utf-8')
```

```
In [22]:  df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 95900 entries, 0 to 95899
Data columns (total 10 columns):
 #   Column         Non-Null Count  Dtype
---  ------         --------------  -----
 0   order_id       95900 non-null  object
 1   cust_id        95900 non-null  object
 2   order_time     95900 non-null  object
 3   order_payment  95900 non-null  float64
 4   pro_id         95900 non-null  object
 5   pro_describe   95900 non-null  object
 6   year           95900 non-null  int64
 7   month          95900 non-null  int64
 8   day            95900 non-null  int64
 9   weekday        95900 non-null  int64
dtypes: float64(1), int64(4), object(5)
memory usage: 7.3+ MB
```
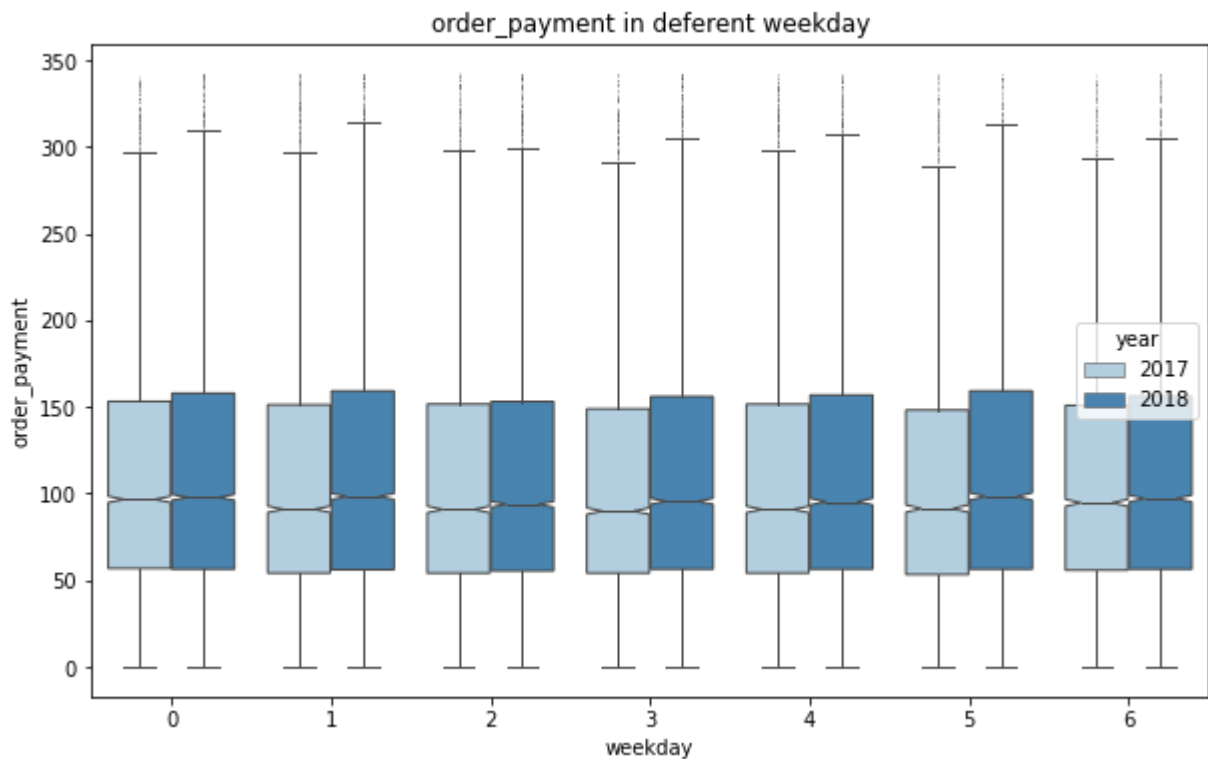
### use seaborn to analyze order_payment in deferent year

```
In [30]:  sns.boxplot(x='year',y='order_payment',data=df,palette='Blues')
          plt.title('order_payment in deferent year')
          plt.show()
```

order_payment in deferent year

**use seaborn to analyze order_payment in diferent weekday**

In [34]:
```python
sns.boxplot(x='weekday',y='order_payment', hue='year',linewidth=1,fliersize=0.05,whis=1
plt.title('order_payment in deferent weekday')
plt.show
```

Out[34]: <function matplotlib.pyplot.show(close=None, block=None)>



order_payment in deferent weekday

**use seaborn to analyze year-by-month total of the customer**

```
In [55]:   # get the different customer of each month
           df_month_customerid = df.groupby(['year','month','cust_id'])['order_payment'].sum().res
           df_month_customerid
```

Out[55]:

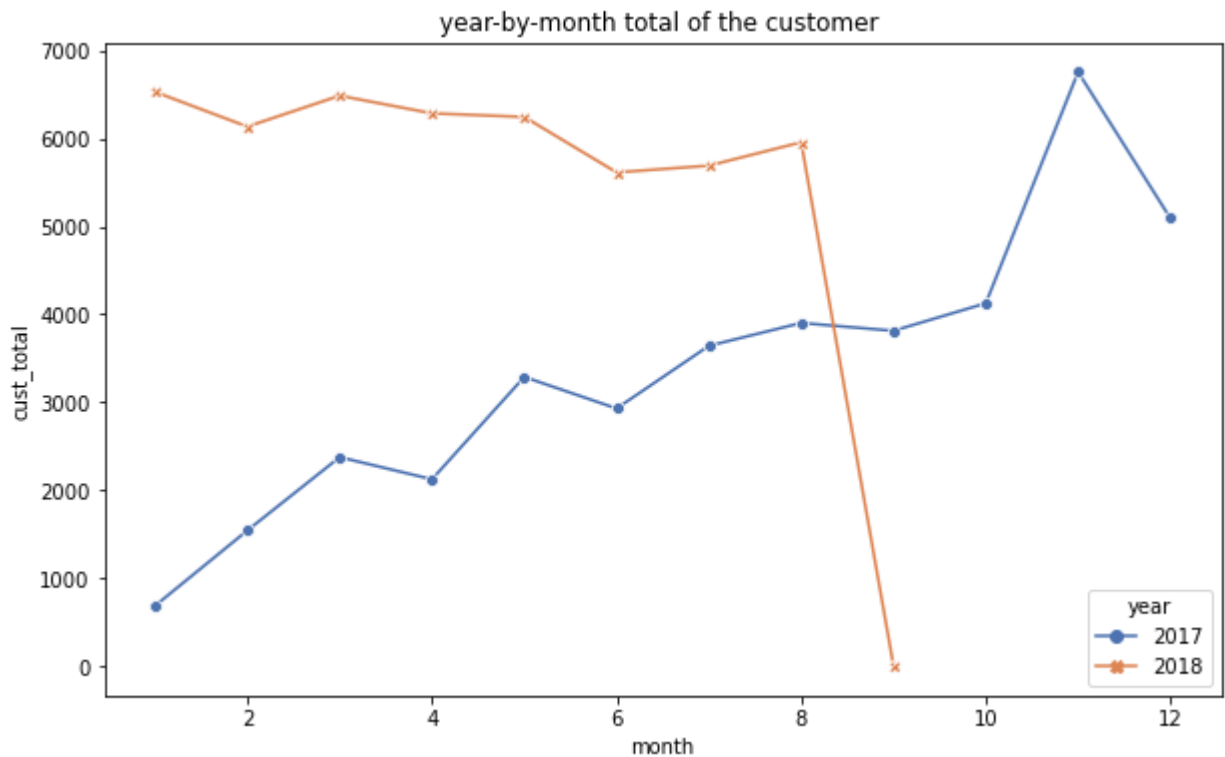| | year | month | cust_id | order_payment |
|---|---|---|---|---|
| **0** | 2017 | 1 | 0040b00970e2139e8c43b647c0da5305 | 41.93 |
| **1** | 2017 | 1 | 0051337a96842850e1ec728dd158f4b3 | 237.99 |
| **2** | 2017 | 1 | 007b7f04a35e02745c23ea706492ca20 | 77.06 |
| **3** | 2017 | 1 | 00f3b3a7cd0b6566435090c7fbda03a2 | 57.51 |
| **4** | 2017 | 1 | 01a0d45a369a4356ac4652584652109a | 45.86 |
| **...** | ... | ... | ... | ... |
| **89222** | 2018 | 8 | ffb3857a7f2f2945434d57e00d0a97a7 | 131.38 |
| **89223** | 2018 | 8 | ffb5eaca500a57b7dd52256fcfc82e12 | 93.63 |
| **89224** | 2018 | 8 | ffe1eab23bff108bf37c973b05d4e9ba | 98.65 |
| **89225** | 2018 | 8 | fff212062d600f2e1d53f3c5d4a25138 | 65.44 |
| **89226** | 2018 | 9 | 4b7decb9b58e2569548b8b4c8e20e8d7 | 166.46 |

89227 rows × 4 columns

```
In [57]:   df_month_customer_all = df_month_customerid.groupby(['year','month'])['cust_id'].count(
           df_month_customer_all
```

Out[57]:

| | year | month | cust_total |
|---|---|---|---|
| **0** | 2017 | 1 | 684 |
| **1** | 2017 | 2 | 1541 |
| **2** | 2017 | 3 | 2375 |
| **3** | 2017 | 4 | 2123 |
| **4** | 2017 | 5 | 3288 |
| **5** | 2017 | 6 | 2929 |
| **6** | 2017 | 7 | 3642 |
| **7** | 2017 | 8 | 3902 |
| **8** | 2017 | 9 | 3813 |
| **9** | 2017 | 10 | 4126 |
| **10** | 2017 | 11 | 6757 |
| **11** | 2017 | 12 | 5093 |
| **12** | 2018 | 1 | 6531 |
| **13** | 2018 | 2 | 6136 |

|    | year | month | cust_total |
|----|------|-------|------------|
| 14 | 2018 | 3     | 6488       |
| 15 | 2018 | 4     | 6287       |
| 16 | 2018 | 5     | 6246       |
| 17 | 2018 | 6     | 5615       |
| 18 | 2018 | 7     | 5692       |
| 19 | 2018 | 8     | 5958       |
| 20 | 2018 | 9     | 1          |

In [63]:
```python
sns.lineplot(x='month',y='cust_total',data=df_month_customer_all,hue='year',style='year
plt.title('year-by-month total of the customer ')
plt.show()
```



year-by-month total of the customer

### use seaborn to analyze the distribution of customer according to the order\'s number

In [69]:
```python
# sum the order number of different customer
df_diff_customer_order = df.groupby('cust_id')['order_id'].count().reset_index().rename
df_diff_customer_order
```

Out[69]:

|   | cust_id | order_total |
|---|---------|-------------|
| 0 | 00012a2ce6f8dcda20d059ce98491703 | 1 |
| 1 | 000161a058600d5901f007fab4c27140 | 1 |
| 2 | 0001fd6190edaaf884bcaf3d49edf079 | 1 |

|  | cust_id | order_total |
|---|---|---|
| **3** | 0002414f95344307404f0ace7a26f1d5 | 1 |
| **4** | 000379cdec625522490c315e70c7a9fb | 1 |
| **...** | ... | ... |
| **89222** | fffcb937e9dd47a13f05ecb8290f4d3e | 1 |
| **89223** | fffecc9f79fd8c764f843e9951b11341 | 3 |
| **89224** | fffeda5b6d849fbd39689bb92087f431 | 1 |
| **89225** | ffff42319e9b2d713724ae527742af25 | 1 |
| **89226** | ffffa3172527f765de70084a7e53aae8 | 1 |

89227 rows × 2 columns

In [78]:
```python
# sum the total of customer
customer_total = len(df_diff_customer_order)
# get the total of customers with different order number
df_customer_diff_order = df_diff_customer_order.groupby('order_total')['cust_id'].count
# get the retio of customers with different order number
ratio_cust__diff_order = (df_customer_diff_order / customer_total*100).round(2).to_fram
ratio_cust__diff_order
```
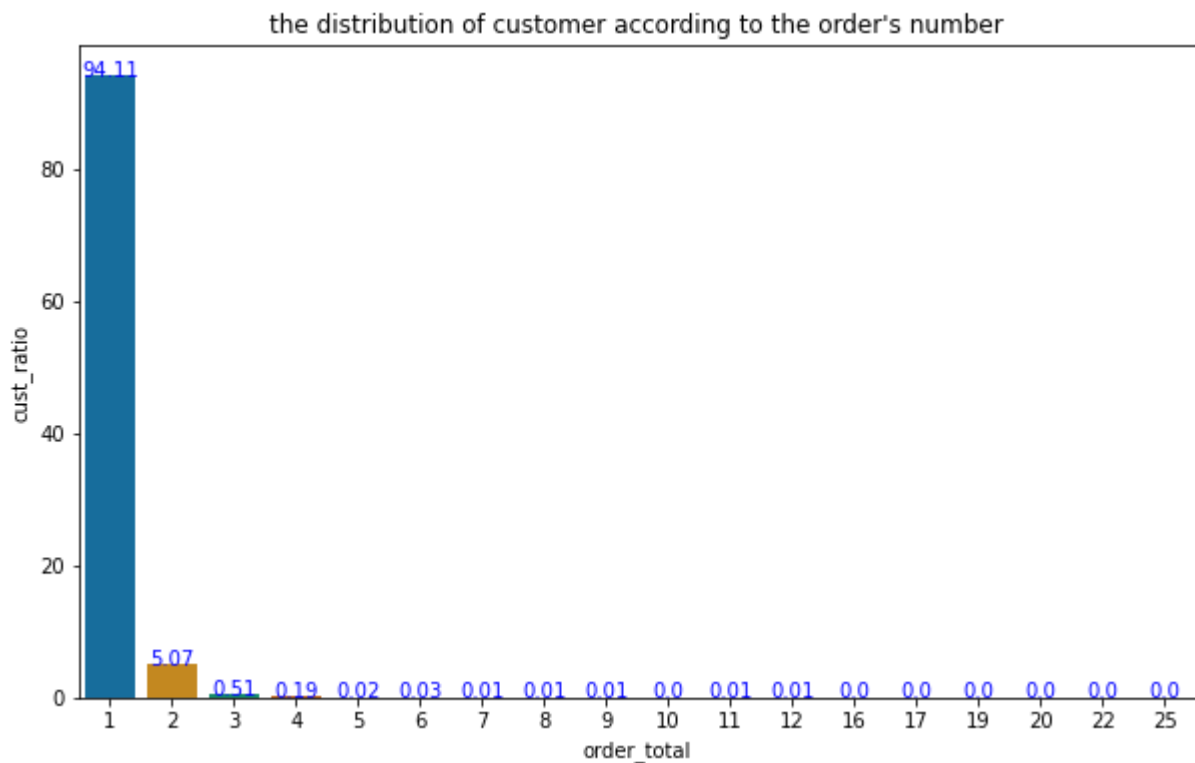
Out[78]:

|  | order_total | cust_ratio |
|---|---|---|
| **0** | 1 | 94.11 |
| **1** | 2 | 5.07 |
| **2** | 3 | 0.51 |
| **3** | 4 | 0.19 |
| **4** | 5 | 0.02 |
| **5** | 6 | 0.03 |
| **6** | 7 | 0.01 |
| **7** | 8 | 0.01 |
| **8** | 9 | 0.01 |
| **9** | 10 | 0.00 |
| **10** | 11 | 0.01 |
| **11** | 12 | 0.01 |
| **12** | 16 | 0.00 |
| **13** | 17 | 0.00 |
| **14** | 19 | 0.00 |
| **15** | 20 | 0.00 |
| **16** | 22 | 0.00 |

|    | order_total | cust_ratio |
|----|-------------|------------|
| **17** | 25 | 0.00 |

```python
# draw the bar chart
g = sns.barplot(x='order_total',y='cust_ratio',data=ratio_cust__diff_order,palette='col
# add the tags for every bar
for index,row in ratio_cust__diff_order.iterrows():
    g.text(row.name,row['cust_ratio'],round(row['cust_ratio'],2),color='blue',ha='cente

plt.xlabel('order_total')
plt.ylabel('cust_ratio')
plt.title('the distribution of customer according to the order\'s number')
plt.show()
```

the distribution of customer according to the order's number