# Data preparation

| | InvoiceNo | StockCode | Description | Quantity | InvoiceDate | UnitPrice | CustomerID | Country |
|---|---|---|---|---|---|---|---|---|
| 0 | 536365 | 85123A | WHITE HANGING HEART T-LIGHT HOLDER | 6 | 2010-12-01 08:26:00 | 2.550 | 17850.000 | United Kingdom |
| 1 | 536365 | 71053 | WHITE METAL LANTERN | 6 | 2010-12-01 08:26:00 | 3.390 | 17850.000 | United Kingdom |
| 2 | 536365 | 84406B | CREAM CUPID HEARTS COAT HANGER | 8 | 2010-12-01 08:26:00 | 2.750 | 17850.000 | United Kingdom |
| 3 | 536365 | 84029G | KNITTED UNION FLAG HOT WATER BOTTLE | 6 | 2010-12-01 08:26:00 | 3.390 | 17850.000 | United Kingdom |
| 4 | 536365 | 84029E | RED WOOLLY HOTTIE WHITE HEART. | 6 | 2010-12-01 08:26:00 | 3.390 | 17850.000 | United Kingdom |

## Missing values

## Feature engeneering

```
RangeIndex: 541909 entries, 0 to 541908
Data columns (total 8 columns):
 #   Column       Non-Null Count    Dtype
---  ------       --------------    -----
 0   InvoiceNo    541909 non-null   object
 1   StockCode    541909 non-null   object
 2   Description  540455 non-null   object
 3   Quantity     541909 non-null   int64
 4   InvoiceDate  541909 non-null   datetime64[ns]
 5   UnitPrice    541909 non-null   float64
 6   CustomerID   406829 non-null   float64
 7   Country      541909 non-null   object
dtypes: datetime64[ns](1), float64(2), int64(1), object(4)
memory usage: 33.1+ MB
```

# Missing values

| | |
|---|---|
| InvoiceNo | 0 |
| StockCode | 0 |
| Description | 1454 |
| Quantity | 0 |
| InvoiceDate | 0 |
| UnitPrice | 0 |
| CustomerID | 135080 |
| Country | 0 |
| Order_price | 0 |
| Order_status | 0 |
| Year | 0 |
| Month | 0 |
| Day | 0 |
| Hour | 0 |
| Minute | 0 |

# Feature engeneering

```python
clust_data["Order_price"]=clust_data["Quantity"]*clust_data["UnitPrice"]
clust_data["Order_status"]=["Done" if order_price>0 else "Cancelled" for order_price in clust_data["Order_price"]]
clust_data["Year"]=clust_data["InvoiceDate"].dt.year
clust_data["Month"]=clust_data["InvoiceDate"].dt.month
clust_data["Day"]=clust_data["InvoiceDate"].dt.day
clust_data["Hour"]=clust_data["InvoiceDate"].dt.hour
clust_data["Minute"]=clust_data["InvoiceDate"].dt.minute
```

```python
clust_data.head(2)
```

| | InvoiceNo | StockCode | Description | Quantity | InvoiceDate | UnitPrice | CustomerID | Country | Order_price | Order_status | Year | Month | Day | Hour | Minute |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 536365 | 85123A | WHITE HANGING HEART T-LIGHT HOLDER | 6 | 2010-12-01 08:26:00 | 2.55 | 17850.0 | United Kingdom | 15.30 | Done | 2010 | 12 | 1 | 8 | 26 |
| 1 | 536365 | 71053 | WHITE METAL LANTERN | 6 | 2010-12-01 08:26:00 | 3.39 | 17850.0 | United Kingdom | 20.34 | Done | 2010 | 12 | 1 | 8 | 26 |

```python
clust_data_=clust_data_.drop(columns=["CustomerID",'StockCode','InvoiceNo','InvoiceDate',"Year"])
```

```python
clust_data_.duplicated().sum()
```
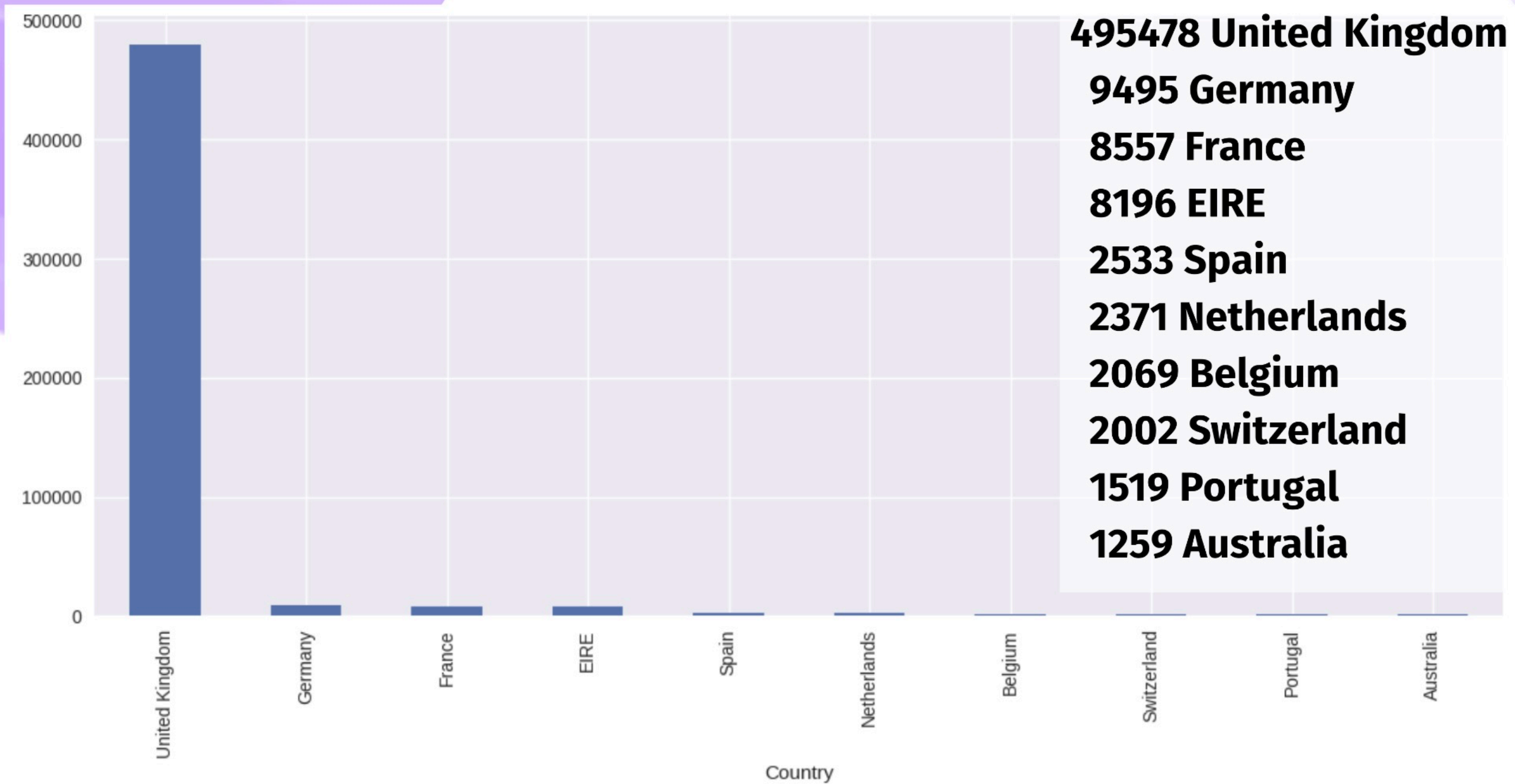
6055

# Correlation

# Data Exploration



495478 United Kingdom
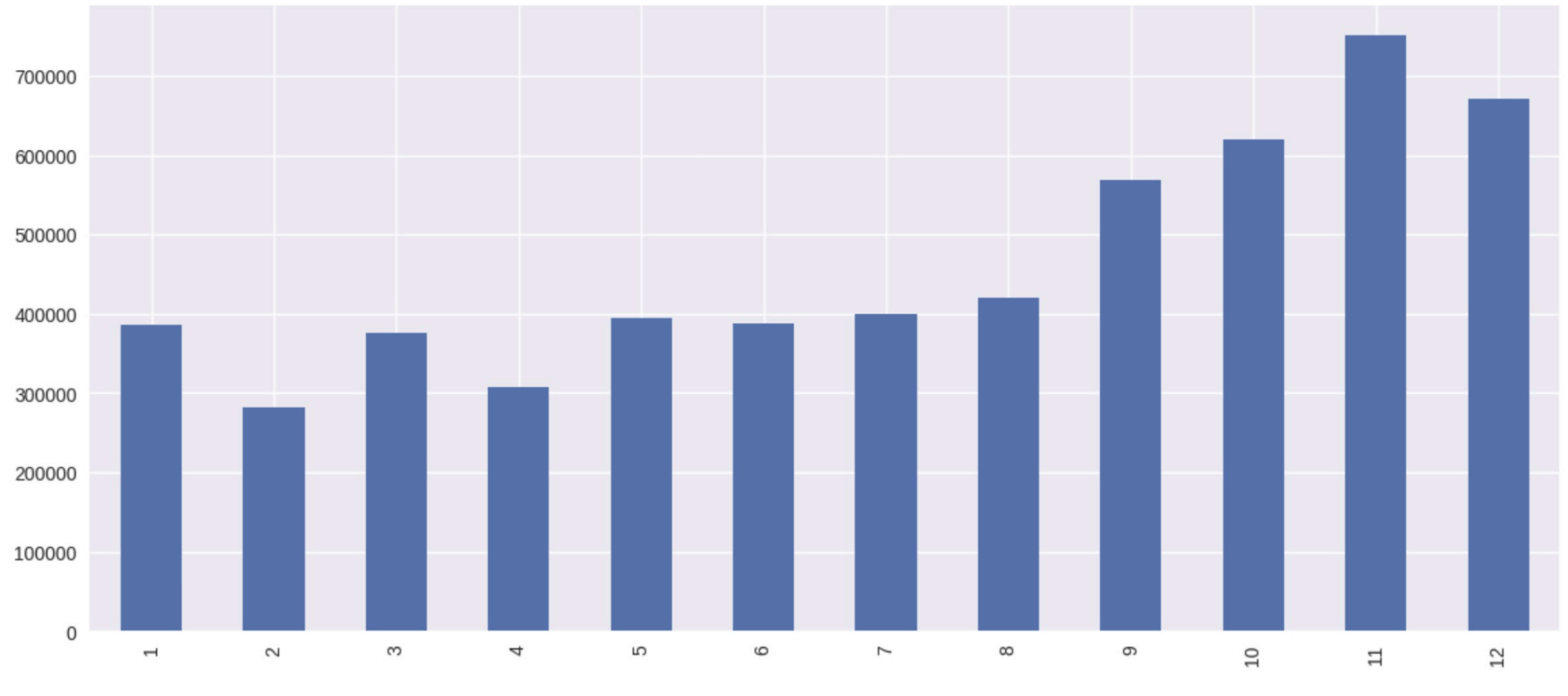9495 Germany
8557 France
8196 EIRE
2533 Spain
2371 Netherlands
2069 Belgium
2002 Switzerland
1519 Portugal
1259 Australia

Quant/ Month

Quant/ Day

Quant/ Hour

Order price

Sales

Quantity

Month

Sales

# Order price

# PCA


Inertia explained per factorial axes
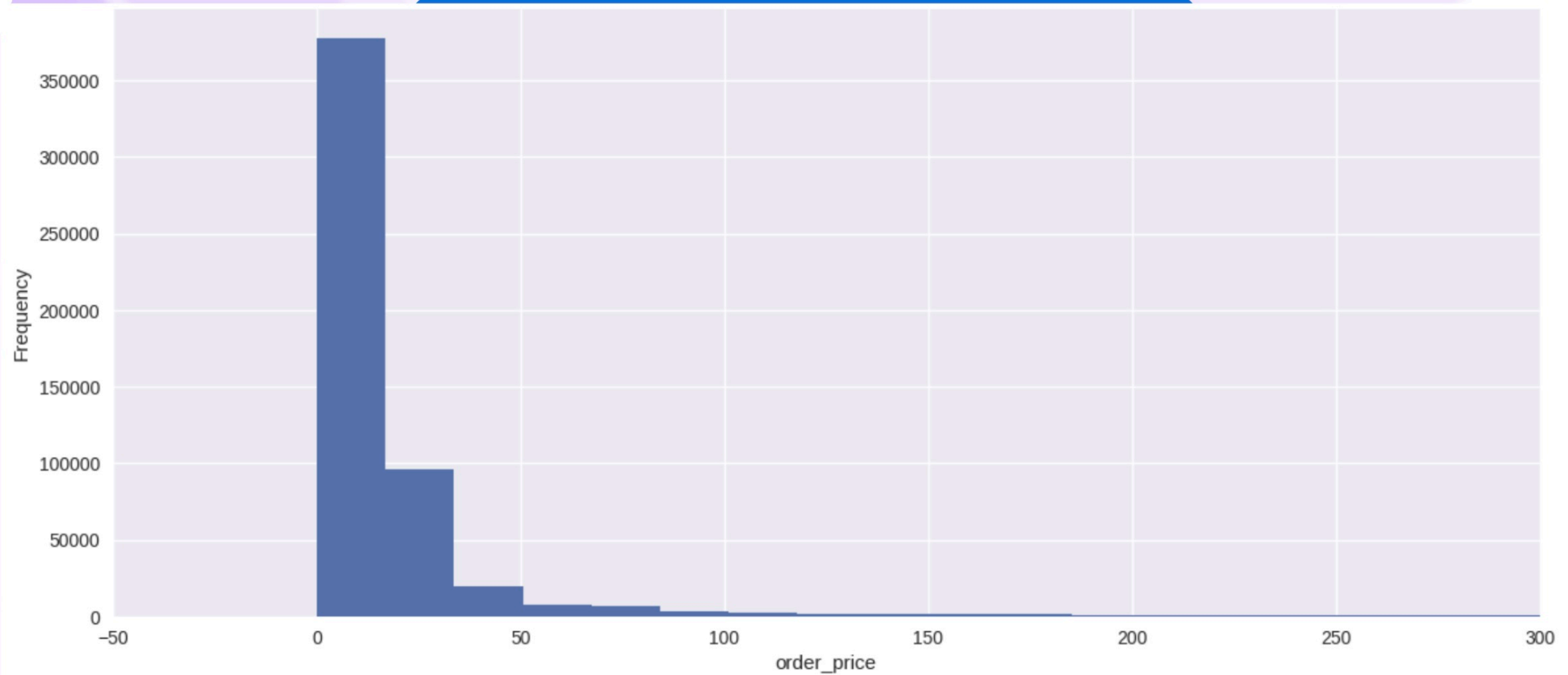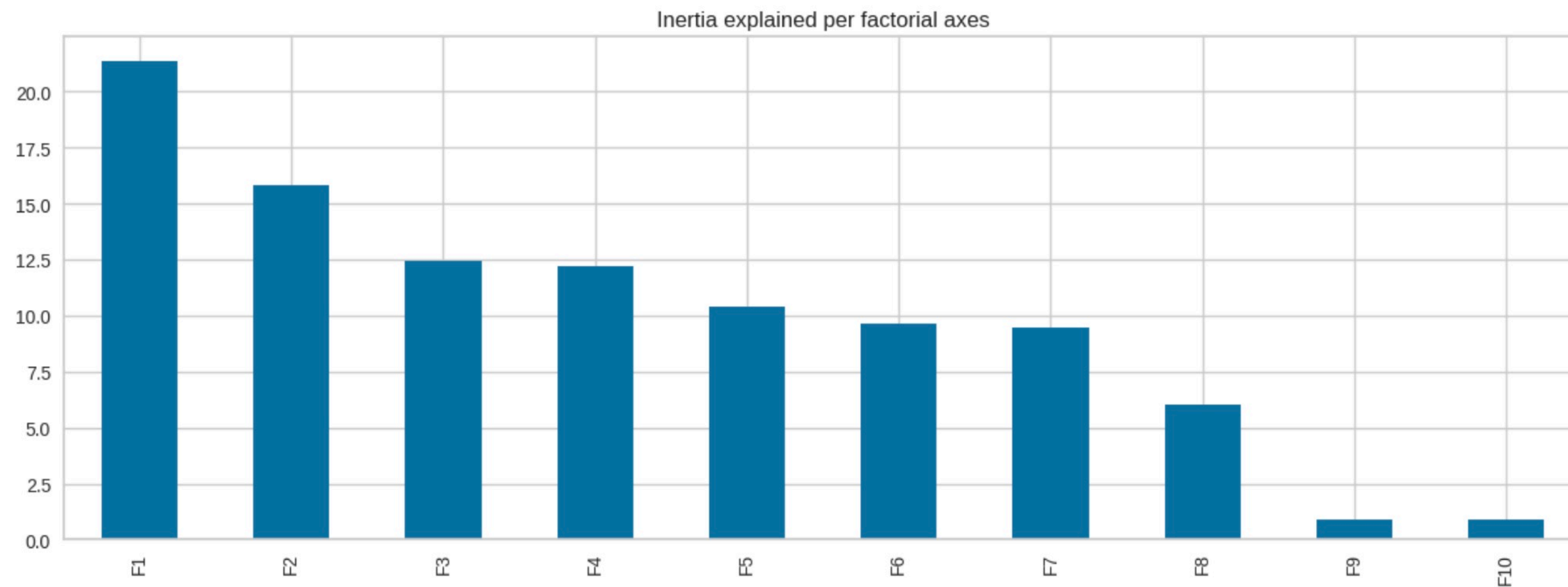
**Choose of K**

**Interclass distance**

**First 8 eigen vector explain 97.4% of the variability**

**From 47 originals columns we keep only 8 columns**

# Elbow plot



Distortion Score Elbow for KMeans Clustering