





```
In [32]: df.duplicated(keep='first').sum()
Out[32]: 0
```

## Treat Outliers

```
In [34]: df.describe(include='all')
Out[34]:
```

	Company	SpecificBeanOrigin	ReviewDate	CocoaPercent	CompanyLocation	Rating	BeanType	BroadBeanOrigin
count	1795	1795	1795.00	1795.00	1795	1795.00	1795	1795
unique	416	1039	nan	nan	60	nan	41	100
top	Soma	Madagascar	nan	nan	USA	nan	Missing	Venezuela
freq	47	57	nan	nan	764	nan	888	214
mean	NaN	NaN	2012.33	0.72	NaN	3.19	NaN	NaN
std	NaN	NaN	2.93	0.06	NaN	0.48	NaN	NaN
min	NaN	NaN	2006.00	0.42	NaN	1.00	NaN	NaN
25%	NaN	NaN	2010.00	0.70	NaN	2.88	NaN	NaN
50%	NaN	NaN	2013.00	0.70	NaN	3.25	NaN	NaN
75%	NaN	NaN	2015.00	0.75	NaN	3.50	NaN	NaN
max	NaN	NaN	2017.00	1.00	NaN	5.00	NaN	NaN

## Type Change

```
In [35]: df.info()
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 1795 entries, 0 to 1794
Data columns (total 8 columns):
# Column Non-Null Count Dtype
---
0 Company 1795 non-null object
1 SpecificBeanOrigin 1795 non-null object
2 ReviewDate 1795 non-null int64
3 CocoaPercent 1795 non-null float64
4 CompanyLocation 1795 non-null object
5 Rating 1795 non-null float64
6 BeanType 1795 non-null object
7 BroadBeanOrigin 1795 non-null object
dtypes: float64(2), int64(1), object(5)
memory usage: 112.3+ KB

In [36]: df['RatingCat'] = df['Rating'].copy()
In [37]: df
Out[37]:
```

	Company	SpecificBeanOrigin	ReviewDate	CocoaPercent	CompanyLocation	Rating	BeanType	BroadBeanOrigin	RatingCat
0	A. Morin	Agua Grande	2016	0.63	France	3.75	Missing	Sao Tome	3.75
1	A. Morin	Kpime	2015	0.70	France	2.75	Missing	Togo	2.75
2	A. Morin	Attane	2015	0.70	France	3.00	Missing	Togo	3.00
3	A. Morin	Akaka	2015	0.70	France	3.50	Missing	Togo	3.50
4	A. Morin	Quilla	2015	0.70	France	3.50	Missing	Peru	3.50
...	...	...	...	...	...	...	...	...	...
1790	Zotter	Peru	2011	0.70	Austria	3.75	Missing	Peru	3.75
1791	Zotter	Congo	2011	0.65	Austria	3.00	Forastero	Congo	3.00
1792	Zotter	Kerala State	2011	0.65	Austria	3.50	Forastero	India	3.50
1793	Zotter	Kerala State	2011	0.62	Austria	3.25	Missing	India	3.25
1794	Zotter	Brazil, Mitzi Blue	2010	0.65	Austria	3.00	Missing	Brazil	3.00

1795 rows x 9 columns

```
In [38]: df[df['RatingCat'] != df['RatingCat'].astype('category')]
In [39]: df.info()
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 1795 entries, 0 to 1794
Data columns (total 9 columns):
# Column Non-Null Count Dtype
---
0 Company 1795 non-null object
1 SpecificBeanOrigin 1795 non-null object
2 ReviewDate 1795 non-null int64
3 CocoaPercent 1795 non-null float64
4 CompanyLocation 1795 non-null object
5 Rating 1795 non-null float64
6 BeanType 1795 non-null object
7 BroadBeanOrigin 1795 non-null object
8 RatingCat 1795 non-null category
dtypes: category(1), float64(2), int64(1), object(5)
memory usage: 114.8+ KB

In [40]: df.describe(include='all')
Out[40]:
```

	Company	SpecificBeanOrigin	ReviewDate	CocoaPercent	CompanyLocation	Rating	BeanType	BroadBeanOrigin	RatingCat
count	1795	1795	1795.00	1795.00	1795	1795.00	1795	1795.00	
unique	416	1039	nan	nan	60	nan	41	100	
top	Soma	Madagascar	nan	nan	USA	nan	Missing	Venezuela	
freq	47	57	nan	nan	764	nan	888	214	
mean	NaN	NaN	2012.33	0.72	NaN	3.19	NaN	NaN	
std	NaN	NaN	2.93	0.06	NaN	0.48	NaN	NaN	
min	NaN	NaN	2006.00	0.42	NaN	1.00	NaN	NaN	
25%	NaN	NaN	2010.00	0.70	NaN	2.88	NaN	NaN	
50%	NaN	NaN	2013.00	0.70	NaN	3.25	NaN	NaN	
75%	NaN	NaN	2015.00	0.75	NaN	3.50	NaN	NaN	
max	NaN	NaN	2017.00	1.00	NaN	5.00	NaN	NaN	

## Hypothesis Testing

The goal of hypothesis testing is to answer the question, "Given a sample and an apparent effect, what is the probability of seeing such an effect by chance?" The first step is to quantify the size of the apparent effect by choosing a test statistic (t-test, ANOVA, etc). The next step is to define a null hypothesis, which is a model of the system based on the assumption that the apparent effect is not real. Then compute the p-value, which is the probability of the null hypothesis being true, and finally interpret the result of the p-value, if the value is low, the effect is said to be statistically significant, which means that the null hypothesis may not be accurate.

### T-Test

We will be using the t-test for independent samples. For the independent t-test, the following assumptions must be met.

- One independent, categorical variable with two levels or group
- One dependent continuous variable
- Independence of the observations. Each subject should belong to only one group. There is no relationship between the observations in each group.
- The dependent variable must follow a normal distribution
- Assumption of homogeneity of variance

State the hypothesis

- $H_0: \mu_1 = \mu_2$  ("there is no difference between cocoa and ratings")
- $H_1: \mu_1 \neq \mu_2$  ("there is a difference between cocoa and ratings")

### T-Test

#### One Sample T-Test

```
In [41]: t, p = scipy.stats.ttest_1samp(a=df.CocoaPercent, popmean=0.72)
In [42]: print(f"T-test value is: ", t)
print(f"p-value value is: ", p)
T-test value is: -2.008291487368794
p-value value is: 0.044761965140545666

In [43]: t, p = scipy.stats.ttest_ind(a=df.CocoaPercent, b=df.RatingCat, equal_var = False)
In [44]: print(f"T-test value is: ", t)
print(f"p-value value is: ", p)
T-test value is: -216.91553551315675
p-value value is: 0.0
There is statistical significance since p-value < 0.05, Null Hypothesis is rejected
```

## Chi-square

State the hypothesis:

- $H_0: \chi^2$  The proportion of Ratings is independent of Bean Type
- $H_1: \chi^2$  The proportion of Ratings who are tenured is associated with Bean Type

```
In [45]: #Create a Cross-tab table
cont_table = pd.crosstab(df['BeanType'], df['RatingCat'])
cont_table
Out[45]:
```

	RatingCat	1.0	1.5	1.75	2.0	2.25	2.5	2.75	3.0	3.25	3.5	3.75	4.0	5.0
BeanType	Amazon	0	0	0	0	0	0	0	0	1	0	0	0	0
Amazon mix	0	0	0	0	0	0	0	0	0	1	0	1	0	0
Amazon, ICS	0	0	0	0	0	0	0	0	0	1	1	0	0	0
Beniano	0	0	0	0	0	0	0	1	0	0	1	1	0	0
Blend	0	0	0	2	0	2	1	9	3	14	4	5	1	0
Blend-Forastero-Criollo	0	0	0	0	0	0	0	0	0	1	0	0	0	0
CCN51	0	0	0	0	0	0	0	0	0	1	0	0	0	0
Criollo	0	0	0	1	1	13	18	39	13	33	21	14	0	0
Criollo (Amarru)	0	0	0	0	0	0	0	1	0	1	0	0	0	0
Criollo (Ocumare 61)	0	0	0	0	0	0	0	0	2	0	0	0	0	0
Criollo (Ocumare 67)	0	0	0	0	0	0	0	0	0	0	0	1	0	0
Criollo (Ocumare 77)	0	0	0	0	0	0	0	0	0	0	1	0	0	0
Criollo (Ocumare)	0	0	0	0	0	0	0	1	0	0	0	0	0	0
Criollo (Porcelana)	0	0	0	0	0	1	0	2	3	1	0	0	0	0
Criollo (Wild)	0	0	0	0	0	0	0	0	0	0	0	1	0	0
Criollo, +	0	0	0	0	0	0	0	0	0	1	0	0	0	0
Criollo, Forastero	0	0	0	0	0	0	0	0	0	1	1	0	0	0
Criollo, Trinitario	0	0	0	0	0	2	4	8	6	13	2	4	0	0
EET	0	0	0	0	0	0	0	0	0	2	1	0	0	0
Forastero	1	1	0	2	1	9	13	18	10	19	9	4	0	0
Forastero (Amelonado)	0	0	0	0	0	0	0	0	0	0	1	0	0	0
Forastero (Arriba)	0	2	0	4	0	8	4	8	2	7	1	1	0	0
Forastero (Arriba) ASS	0	0	0	1	0	1	0	2	1	1	0	0	0	0
Forastero (Arriba) ASS5	0	0	0	0	0	0	0	0	0	1	0	0	0	0
Forastero (Catongo)	0	0	0	0	0	0	0	0	1	1	0	0	0	0
Forastero (Nacional)	0	0	0	1	0	2	6	13	5	13	9	3	0	0
Forastero (Parazinho)	0	0	0	0	0	0	1	0	0	4	2	1	0	0
Forastero(Arriba, CCN)	0	0	0	0	0	0	0	1	0	0	0	0	0	0
Forastero, Trinitario	0	0	0	0	0	0	0	1	0	0	0	0	0	0
Matina	0	0	0	0	0	0	0	1	0	1	1	0	0	0
Missing	3	7	2	19	9	67	150	146	182	169	101	33	0	0
Nacional	0	0	0	0	0	1	0	0	0	1	0	0	0	0
Nacional (Arriba)	0	0	0	0	0	0	1	0	1	0	1	0	0	0
Trinitario	0	0	1	2	2	19	61	85	72	103	48	25	1	0
Trinitario (85% Criollo)	0	0	0	0	0	0	0	0	0	0	1	1	0	0
Trinitario (Amelonado)	0	0	0	0	0	0	0	0	1	0	0	0	0	0
Trinitario (Scavina)	0	0	0	0	0	0	0	0	0	1	0	0	0	0
Trinitario, Criollo	0	0	0	1	2	0	2	1	2	1	1	0	0	0
Trinitario, Forastero	0	0	0	0	0	0	0	2	0	0	0	0	0	0
Trinitario, Nacional	0	0	0	0	0	0	0	0	0	0	1	0	0	0
Trinitario, TCGA	0	0	0	0	0	0	0	0	0	0	1	0	0	0

```
In [46]: chi_square = scipy.stats.chi2_contingency(cont_table, correction = True)
In [47]: print(f"Chi score is", chi_square[0])
Chi score is 410.97839899393007
In [48]: print(f"P-value is", chi_square[1])
P-value is 0.9898999917428333
In [49]: print(f"Degrees of freedom is", chi_square[2])
Degrees of freedom is 480
P-value > 0.05 hence both variables are independent of each other
```

### Correlation

State the hypothesis:

- $H_0: \rho$  Cocoa Percent is not correlated with Rating
- $H_1: \rho$  Cocoa Percent is correlated with Rating

```
In [50]: pearson_correlation = scipy.stats.pearsonr(df['CocoaPercent'], df['Rating'])
In [51]: print(f"Pearson's correlation coefficient is", pearson_correlation[0])
Pearson's correlation coefficient is -0.16503542231877807
In [52]: print(f"P-value is", pearson_correlation[1])
P-value is 1.984730626881529e-12
P-value < 0.05 hence Cocoa Percent has impact on Ratings. Reject Null Hypothesis
```

## Select Chocolate Ratings >= 4.0 and Cocoa Percent >= 0.8

```
In [53]: df.head()
Out[53]:
```

	Company	SpecificBeanOrigin	ReviewDate	CocoaPercent	CompanyLocation	Rating	BeanType	BroadBeanOrigin	RatingCat
0	A. Morin	Agua Grande	2016	0.63	France	3.75	Missing	Sao Tome	3.75
1	A. Morin	Kpime	2015	0.70	France	2.75	Missing	Togo	2.75
2	A. Morin	Attane	2015	0.70	France	3.00	Missing	Togo	3.00
3	A. Morin	Akaka	2015	0.70	France	3.50	Missing	Togo	3.50
4	A. Morin	Quilla	2015	0.70	France	3.50	Missing	Peru	3.50

```
In [54]: df2 = df[df['Rating'] >= 4.0]
In [55]: df2.head()
Out[55]:
```

	Company	SpecificBeanOrigin	ReviewDate	CocoaPercent	CompanyLocation	Rating	BeanType	BroadBeanOrigin	RatingCat
9	A. Morin	Pabino	2014	0.70	France	4.00	Missing	Peru	4.00
17	A. Morin	Chuao	2013	0.70	France	4.00	Trinitario	Venezuela	4.00
20	A. Morin	Chanchamayo Province	2013	0.63	France	4.00	Missing	Peru	4.00
54	Amano	Morobe	2011	0.70	USA	4.00	Missing	Papua New Guinea	4.00
56	Amano	Guayas	2010	0.70	USA	4.00	Missing	Ecuador	4.00

```
In [56]: df2.CocoaPercent.value_counts()
Out[56]:
```

CocoaPercent	count
0.70	45
0.75	17
0.72	11
0.74	4
0.64	4
0.63	3
0.67	2
0.68	2
0.69	2
0.65	2
0.80	1
0.60	1
0.78	1
0.71	1
0.66	1
0.73	1
0.88	1

```
In [57]: fig = plt.figure(figsize=(20,40))
plt.subplot(7,2,1)
plt.title("Ratings Counts")
sns.countplot(df2.Rating)

plt.subplot(7,2,2)
plt.title("Year Counts")
sns.countplot(df2.ReviewDate)

plt.subplot(7,2,3)
plt.title("BarChart")
sns.barplot(x=df2.ReviewDate, y=df2.Rating, data=df2)

plt.subplot(7,2,4)
plt.title("BarChart")
sns.barplot(x=df2.BroadBeanOrigin, y=df2.Rating, data=df2)

plt.subplot(7,2,5)
plt.title("BarChart")
sns.barplot(x=df2.BeanType, y=df2.Rating, data=df2)

plt.subplot(7,2,6)
plt.title("BarChart")
sns.barplot(x=df2.CompanyLocation, y=df2.Rating, data=df2)

plt.subplot(7,2,7)
plt.title("BarChart")
sns.barplot(x=df2.SpecificBeanOrigin, y=df2.Rating, data=df2)

plt.subplot(7,2,8)
plt.title("BarChart")
sns.barplot(x=df2.Company, y=df2.Rating, data=df2)

plt.subplot(7,2,9)
plt.title("Scatterplot")
sns.scatterplot(x=df2.CocoaPercent, y=df2.Rating, data=df2)

plt.subplot(7,2,10)
plt.title("")
sns.scatterplot()
```

