

EDA: Exploratory Data Analysis

EDA + DPP

Data Preprocessing

EDA Checklist

- Table-level analysis
- Per-column analysis and processing
- Cross-column analysis and processing
- Cross-table analysis

EDA: Table-Level Analysis

```
df.describe()
```

```
!df.info(verbose=True, memory_usage=True)  
!df.describe(include='all')
```

	Company	Age(Range of Years)	Specialties (Main Programs)	Revenue Per Month	ROI	Revenue/Unit/Sec	Success/Percent	Company/Innovation	Rating	Rank/Type	Overall Score/Design
count	1245	1245	1245	1245	1245	1245	1245	1245	1245	1245	1245
unique	415	1033	NA	NA	NA	NA	40	93	NA	41	133
top	Google	Open	Machine Learning	NA	NA	NA	75%	USA	NA	NA	Veranda
freq	47	57	NA	NA	NA	NA	872	754	NA	887	214
mean	NA	NA	NA	1235.434735	2012.305048	NA	NA	NA	2.165022	NA	NA
std	NA	NA	NA	527.892865	5.007910	NA	NA	NA	0.570000	NA	NA
min	NA	NA	NA	4.000000	1988.000000	NA	NA	NA	1.000000	NA	NA
25%	NA	NA	NA	570.000000	2010.000000	NA	NA	NA	1.875000	NA	NA
50%	NA	NA	NA	1239.000000	2013.000000	NA	NA	NA	2.250000	NA	NA
75%	NA	NA	NA	532.000000	2015.000000	NA	NA	NA	2.500000	NA	NA
max	NA	NA	NA	457.000000	2017.000000	NA	NA	NA	3.000000	NA	NA

EDA: Table-Level Analysis

```
df.head()      # First few rows of data
df.tail()      # Last few rows of data
df.sample(10)  # Random sample of data
df.columns     # Column labels
df.dtypes      # Data types of columns
```

EDA: Per Column Analysis

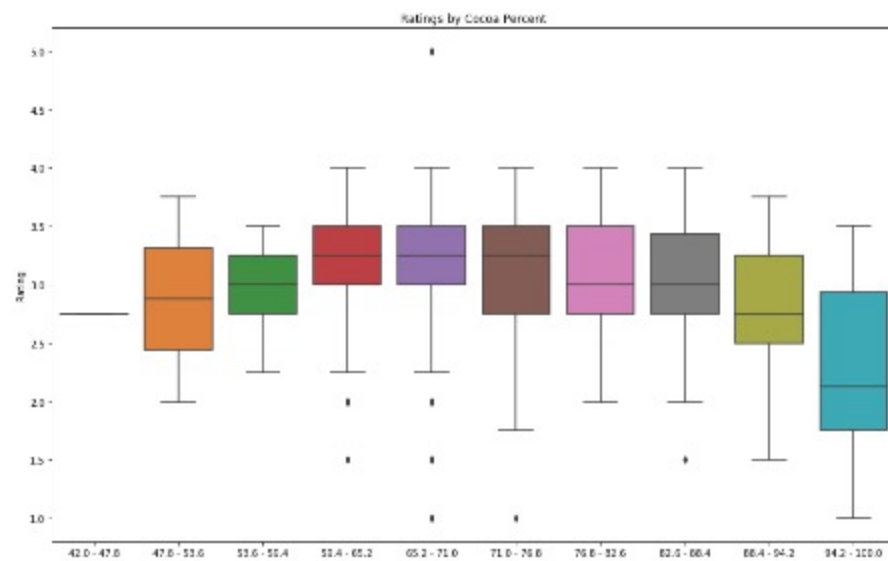
```
df[column].isna().any() # Missing values?  
df[column].describe()  # Statistics for column  
df[column].hist()       # Histogram  
df[column].sample(10)   # Data sample
```

EDA: Cross-Column Analysis

```
1 df.corr()
```

	Cocoa Percent	REF	Review Date	Rating
Cocoa Percent	1.000000	0.037791	0.038177	-0.164820
REF	0.037791	1.000000	0.985076	0.101030
Review Date	0.038177	0.985076	1.000000	0.099846
Rating	-0.164820	0.101030	0.099846	1.000000

EDA: Cross-Column Analysis



EDA: Cross-Table Analysis

	Country	Latitude	Longitude
count	58	58.000000	58.00000
unique	58	NaN	NaN
top	South Korea	NaN	NaN
freq	1	NaN	NaN
mean	NaN	23.392586	-0.45931
std	NaN	29.271631	71.95592
min	NaN	-41.190000	-99.10000
25%	NaN	5.387500	-64.97250
50%	NaN	23.785000	1.07000
75%	NaN	50.395000	20.51250
max	NaN	64.100000	178.30000

EDA: Data Preprocessing

- Missing values
- Invisible values
- Placeholder values
- Values out of the expected range
- Incorrect data types
- Inconsistent values (USA vs. U.S.A.)

EDA: *Not Necessarily* Data Preprocessing

- Normalization
- Standardization
- One-hot encoding

Preprocessing Strategies

- Fill in missing data
- Fix broken data
- Remove problematic columns
- Remove problematic rows

Your Turn!

Chocolate Bar Ratings

- Lab Part 1: Individual tables and columns
- Lab Part 2: Cross tables and columns