

2.2.1 Using Comments in Python

```
'''
This is a multiline comment
'''
#This is a single line comment
```

```
Out[1]: 'InThis is a multiline commentIn'
```

2.2.3 Importing Packages in Python

```
#Importing packages
import pandas as pd
import numpy as np
from sklearn.tree import DecisionTreeClassifier, export_graphviz
import scipy.stats
import matplotlib.pyplot as plt
import seaborn as sns
```

2.2.4 Getting Data into Python

```
#Importing the bank marketing trained dataset
bank_train = pd.read_csv('C:/School/DS-C-530/DataSets/bank_marketing_training')
```

```
#Displaying the data
bank_train
```

| | age | job | marital | education | default | housing | loan | contact | month | day.of.week | ... | campaign | days.since.previous |
|-------|-----|-------------|---------|---------------------|---------|---------|------|-----------|-------|-------------|-----|----------|---------------------|
| 0 | 56 | housemaid | married | basic.4y | no | no | no | telephone | may | mon | ... | 1 | 999 |
| 1 | 57 | services | married | high.school | unknown | no | no | telephone | may | mon | ... | 1 | 999 |
| 2 | 41 | blue-collar | married | unknown | unknown | no | no | telephone | may | mon | ... | 1 | 999 |
| 3 | 25 | services | single | high.school | no | yes | no | telephone | may | mon | ... | 1 | 999 |
| 4 | 29 | blue-collar | single | high.school | no | no | yes | telephone | may | mon | ... | 1 | 999 |
| ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... |
| 26869 | 36 | admin. | married | university.degree | no | no | no | cellular | nov | fri | ... | 2 | 999 |
| 26870 | 37 | admin. | married | university.degree | no | yes | no | cellular | nov | fri | ... | 1 | 999 |
| 26871 | 29 | unemployed | single | basic.4y | no | yes | no | cellular | nov | fri | ... | 1 | 999 |
| 26872 | 73 | retired | married | professional.course | no | yes | no | cellular | nov | fri | ... | 1 | 999 |
| 26873 | 46 | blue-collar | married | professional.course | no | no | no | cellular | nov | fri | ... | 1 | 999 |

26874 rows × 21 columns

2.2.5 Saving Output in Python

```
#Creating a contingency table with no saving
pd.crosstab(bank_train['previous_outcome'], bank_train['response'])
```

| | response | no | yes |
|------------------|-------------|-------|------|
| previous_outcome | failure | 2390 | 385 |
| | nonexistent | 21176 | 2034 |
| | success | 320 | 569 |

```
#Making a contingency table with with saving
crosstab_01 = pd.crosstab(bank_train['previous_outcome'], bank_train['response'])
```

```
#Display the data
crosstab_01
```

| | response | no | yes |
|------------------|-------------|-------|------|
| previous_outcome | failure | 2390 | 385 |
| | nonexistent | 21176 | 2034 |
| | success | 320 | 569 |

2.2.6 Accessing Records and Variables in Python

```
#Retrieving one row
bank_train.loc[0]
```

```
Out[9]: age      56
job      housemaid
marital   married
education basic.4y
default    no
housing    no
loan       no
contact    telephone
month      may
day.of.week mon
duration    261
campaign    999
days.since.previous 0
previous_outcome nonexistent
emp.var.rate      1.1
cons.price.idx    93.994
cons.conf.idx     -36.4
euribor3m         4.457
nr.employed      5191
response         no
dtype: object
```

```
#Extracting three rows
bank_train.loc[0:2,3:11]
```

| | age | job | marital | education | default | housing | loan | contact | month | day.of.week | ... | campaign | days.since.previous | previous |
|---|-----|-------------|---------|-------------|---------|---------|------|-----------|-------|-------------|-----|----------|---------------------|----------|
| 0 | 56 | housemaid | married | basic.4y | no | no | no | telephone | may | mon | ... | 1 | 999 | C |
| 1 | 57 | services | married | high.school | unknown | no | no | telephone | may | mon | ... | 1 | 999 | C |
| 2 | 41 | blue-collar | married | unknown | unknown | no | no | telephone | may | mon | ... | 1 | 999 | C |
| 3 | 25 | services | single | high.school | no | yes | no | telephone | may | mon | ... | 1 | 999 | C |

3 rows × 21 columns

```
#Extracting 10 rows using different example
bank_train.loc[0:10]
```

| | age | job | marital | education | default | housing | loan | contact | month | day.of.week | ... | campaign | days.since.previous | prev |
|----|-----|-------------|----------|-------------|---------|---------|------|-----------|-------|-------------|-----|----------|---------------------|------|
| 0 | 56 | housemaid | married | basic.4y | no | no | no | telephone | may | mon | ... | 1 | 999 | |
| 1 | 57 | services | married | high.school | unknown | no | no | telephone | may | mon | ... | 1 | 999 | |
| 2 | 41 | blue-collar | married | unknown | unknown | no | no | telephone | may | mon | ... | 1 | 999 | |
| 3 | 25 | services | single | high.school | no | yes | no | telephone | may | mon | ... | 1 | 999 | |
| 4 | 29 | blue-collar | single | high.school | no | no | yes | telephone | may | mon | ... | 1 | 999 | |
| 5 | 57 | housemaid | divorced | basic.4y | no | yes | no | telephone | may | mon | ... | 1 | 999 | |
| 6 | 35 | blue-collar | married | basic.6y | no | no | yes | telephone | may | mon | ... | 1 | 999 | |
| 7 | 39 | management | single | basic.9y | unknown | no | no | telephone | may | mon | ... | 1 | 999 | |
| 8 | 30 | unemployed | married | high.school | no | no | no | telephone | may | mon | ... | 1 | 999 | |
| 9 | 55 | retired | single | high.school | no | yes | no | telephone | may | mon | ... | 1 | 999 | |
| 10 | 41 | technician | single | high.school | no | yes | no | telephone | may | mon | ... | 1 | 999 | |

11 rows × 21 columns

```
#Extracting by column
bank_train['age']
```

```
Out[12]: 0      56
1      57
2      41
3      25
4      29
```

```
In [13]: #Extracting more columns
bank_train[['age', 'job']]
```

```
Out[13]:
```

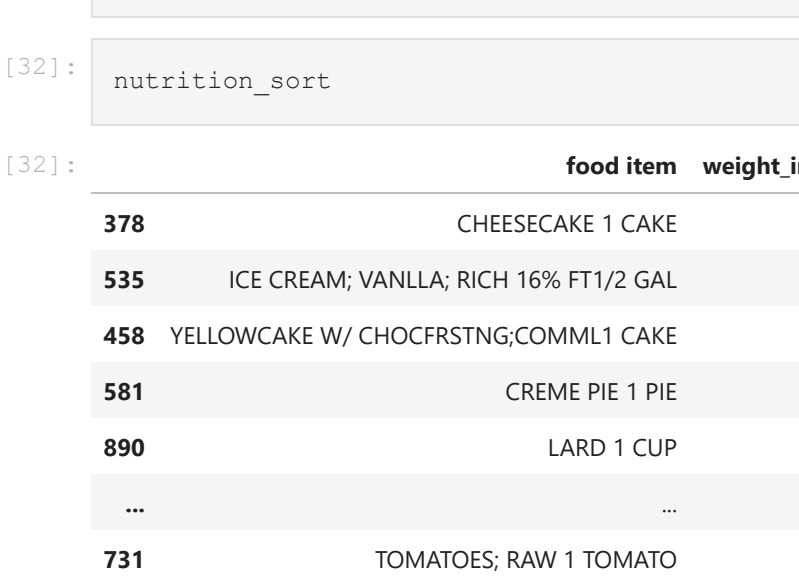
| | age | job |
|-------|-----|-------------|
| 0 | 56 | housemaid |
| 1 | 57 | services |
| 2 | 41 | blue-collar |
| 3 | 25 | services |
| 4 | 29 | blue-collar |
| ... | ... | ... |
| 26869 | 36 | admin. |
| 26870 | 37 | admin |
| 26871 | 29 | unemployed |
| 26872 | 73 | retired |
| 26873 | 46 | blue-collar |

26874 rows × 2 columns

2.2.7 Setting Up Graphics in Python

```
bank_train['age'].plot(kind='hist',title = 'Histogram of Age')
```

```
Out[14]: <AxesSubplot: title='center': 'Histogram of Age', ylabel='Frequency'>
```



Part 2: Data Preparation in Python

```
#Importing the data set
nutrition = pd.read_csv('C:/School/DS-C-530/DataSets/nutrition_subset')
```

```
#Display data
nutrition
```

| | | food item | weight_in_grams | saturated_fat | cholesterol | |
|-----|-----|---------------------------------------|-----------------|---------------|-------------|-----|
| 0 | 1 | GELATIN; DRY 1 ENVELP | 7.00 | 0.0 | 0 | 0 |
| 1 | 2 | SEAWEED; SPIRULINA; DRIED 1 OZ | 28.35 | 0.8 | 0 | 0 |
| 2 | 3 | YEAST; BAKERS; DRY; ACTIVE 1 PKG | 7.00 | 0.0 | 0 | 0 |
| 3 | 4 | PARMESAN CHEESE; GRATED 1 PKG | 28.35 | 5.4 | 22 | 3 |
| 4 | 5 | PARMESAN CHEESE; GRATED 1 CUP | 100.00 | 19.1 | 79 | 4 |
| ... | ... | ... | ... | ... | ... | ... |
| 956 | 957 | COFFEE; BREWED 6 FL OZ | 180.00 | 0.0 | 0 | 956 |
| 957 | 958 | TEA; BREWED 8 FL OZ | 240.00 | 0.0 | 0 | 957 |
| 958 | 959 | TEA; INSTANT;PREP;DRUNSWEETEND8 FL OZ | 241.00 | 0.0 | 0 | 958 |
| 959 | 960 | LETTUCE; BUTTERHEAD; RAW;LEAVE1 LEAF | 15.00 | 0.0 | 0 | 959 |
| 960 | | SALT 1 TSP | 5.50 | 0.0 | 0 | 960 |

961 rows × 4 columns

1. Sort the data set by the saturated fat variable

```
#Adding an index
nutrition['index'] = pd.Series(range(0,961))
```

```
#Display data
nutrition
```

| | | food item | weight_in_grams | saturated_fat | cholesterol | index |
|-----|-----|---------------------------------------|-----------------|---------------|-------------|-------|
| 0 | 1 | GELATIN; DRY 1 ENVELP | 7.00 | 0.0 | 0 | 0 |
| 1 | 2 | SEAWEED; SPIRULINA; DRIED 1 OZ | 28.35 | 0.8 | 0 | 1 |
| 2 | 3 | YEAST; BAKERS; DRY; ACTIVE 1 PKG | 7.00 | 0.0 | 0 | 2 |
| 3 | 4 | PARMESAN CHEESE; GRATED 1 OZ | 28.35 | 5.4 | 22 | 3 |
| 4 | 5 | PARMESAN CHEESE; GRATED 1 CUP | 100.00 | 19.1 | 79 | 4 |
| ... | ... | ... | ... | ... | ... | ... |
| 956 | 957 | COFFEE; BREWED 6 FL OZ | 180.00 | 0.0 | 0 | 956 |
| 957 | 958 | TEA; BREWED 8 FL OZ | 240.00 | 0.0 | 0 | 957 |
| 958 | 959 | TEA; INSTANT;PREP;DRUNSWEETEND8 FL OZ | 241.00 | 0.0 | 0 | 958 |
| 959 | 960 | LETTUCE; BUTTERHEAD; RAW;LEAVE1 LEAF | 15.00 | 0.0 | 0 | 959 |
| 960 | | SALT 1 TSP | 5.50 | 0.0 | 0 | 960 |

961 rows × 5 columns

```
nutrition_sort = nutrition.sort_values(['saturated_fat'], ascending=False)
```

```
#Display data
nutrition_sort
```

| | | food item | weight_in_grams | saturated_fat | cholesterol | index |
|-----|-----|--|-----------------|---------------|-------------|-------|
| 378 | 535 | CHEESECAKE 1 CAKE | 1110.0 | 118.9 | 2053 | 378 |
| 535 | 458 | ICE CREAM; VANILLA; RICH 16% FT1/2 GAL | 1188.0 | 118.3 | 703 | 535 |
| 458 | 458 | YELLOWCAKE W/ CHOCFRTSTNG;COMM;1 CAKE | 1108.0 | 92.0 | 609 | 458 |
| 581 | 581 | CREME PIE 1 PIE | 910.0 | 90.1 | 46 | 581 |
| 890 | 890 | LARD 1 CUP | 205.0 | 80.4 | 195 | 890 |
| ... | ... | ... | ... | ... | ... | ... |
| 731 | 730 | TOMATOES; RAW 1 TOMATO | 123.0 | 0.0 | 0 | 731 |
| 730 | 727 | RASPBERRIES; RAW 1 CUP | 123.0 | 0.0 | 0 | 730 |
| 727 | 726 | TOMATO JUICE; CANNED W/O SALT 1 CUP | 244.0 | 0.0 | 0 | 727 |
| 726 | 960 | TOMATO JUICE; CANNED WITH SALT1 CUP | 244.0 | 0.0 | 0 | 726 |
| 960 | | SALT 1 TSP | 5.5 | 0.0 | 0 | 960 |

961 rows × 5 columns

```
nutrition_sort.head(5)
```

| | | food item | weight_in_grams | saturated_fat | cholesterol | index |
|-----|-----|--|-----------------|---------------|-------------|-------|
| 378 | 535 | CHEESECAKE 1 CAKE | 1110.0 | 118.9 | 2053 | 378 |
| 535 | 458 | ICE CREAM; VANILLA; RICH 16% FT1/2 GAL | 1188.0 | 118.3 | 703 | 535 |
| 458 | 458 | YELLOWCAKE W/ CHOCFRTSTNG;COMM;1 CAKE | 1108.0 | 92.0 | 609 | 458 |
| 581 | 581 | CREME PIE 1 PIE | 910.0 | 90.1 | 46 | 581 |
| 890 | 890 | LARD 1 CUP | 205.0 | 80.4 | 195 | 890 |
| ... | ... | ... | ... | ... | ... | ... |
| 731 | 730 | TOMATOES; RAW 1 TOMATO | 123.0 | 0.0 | 0 | 731 |
| 730 | 727 | RASPBERRIES; RAW 1 CUP | 123.0 | 0.0 | 0 | 730 |
| 727 | 726 | TOMATO JUICE; CANNED W/O SALT 1 CUP | 244.0 | 0.0 | 0 | 727 |
| 726 | 960 | TOMATO JUICE; CANNED WITH SALT1 CUP | 244.0 | 0.0 | 0 | 726 |
| 960 | | SALT 1 TSP | 5.5 | 0.0 | 0 | 960 |

961 rows × 5 columns

```
nutrition_sort.head(5)
```

| | | food item | weight_in_grams | saturated_fat | cholesterol | index |
|-----|-----|--|-----------------|---------------|-------------|-------|
| 378 | 535 | CHEESECAKE 1 CAKE | 1110.0 | 118.9 | 2053 | 378 |
| 535 | 458 | ICE CREAM; VANILLA; RICH 16% FT1/2 GAL | 1188.0 | 118.3 | 703 | 535 |
| 458 | 458 | YELLOWCAKE W/ CHOCFRTSTNG;COMM;1 CAKE | 1108.0 | 92.0 | 609 | 458 |
| 581 | 581 | CREME PIE 1 PIE | 910.0 | 90.1 | 46 | 581 |
| 890 | 890 | LARD 1 CUP | 205.0 | 80.4 | 195 | 890 |
| ... | ... | ... | ... | ... | ... | ... |
| 731 | 730 | TOMATOES; RAW 1 TOMATO | 123.0 | 0.0 | 0 | 731 |
| 730 | 727 | RASPBERRIES; RAW 1 CUP | 123.0 | 0.0 | 0 | 730 |
| 727 | 726 | TOMATO JUICE; CANNED W/O SALT 1 CUP | 244.0 | 0.0 | 0 | 727 |
| 726 | 960 | TOMATO JUICE; CANNED WITH SALT1 CUP | 244.0 | 0.0 | 0 | 726 |
| 960 | | SALT 1 TSP | 5.5 | 0.0 | 0 | 960 |

961 rows × 5 columns

1. Derive a new variable, saturated_fat_per_gram

```
nutrition['saturated_fat_per_gram'] = nutrition['saturated_fat'].div(nutrition['weight_in_grams'].values)
```

```
#Display data
nutrition
```

| | | food item | weight_in_grams | saturated_fat | cholesterol | index | saturated_fat_per_gram |
|-----|-----|---------------------------------------|-----------------|---------------|-------------|-------|------------------------|
| 0 | 1 | GELATIN; DRY 1 ENVELP | 7.00 | 0.0 | 0 | 0 | 0.000000 |
| 1 | 2 | SEAWEED; SPIRULINA; DRIED 1 OZ | 28.35 | 0.8 | 0 | 1 | 0.028219 |
| 2 | 3 | YEAST; BAKERS; DRY; ACTIVE 1 PKG | 7.00 | 0.0 | 0 | 2 | 0.000000 |
| 3 | 4 | PARMESAN CHEESE; GRATED 1 OZ | 28.35 | 5.4 | 22 | 3 | 0.190476 |
| 4 | 5 | PARMESAN CHEESE; GRATED 1 CUP | 100.00 | 19.1 | 79 | 4 | 0.191000 |
| ... | ... | ... | ... | ... | ... | ... | ... |
| 956 | 957 | COFFEE; BREWED 6 FL OZ | 180.00 | 0.0 | 0 | 956 | 0.000000 |
| 957 | 958 | TEA; BREWED 8 FL OZ | 240.00 | 0.0 | 0 | 957 | 0.000000 |
| 958 | 959 | TEA; INSTANT;PREP;DRUNSWEETEND8 FL OZ | 241.00 | 0.0 | 0 | 958 | 0.000000 |
| 959 | 960 | LETTUCE; BUTTERHEAD; RAW;LEAVE1 LEAF | 15.00 | 0.0 | 0 | 959 | 0.000000 |
| 960 | | SALT 1 TSP | 5.50 | 0.0 | 0 | 960 | 0.000000 |

961 rows × 6 columns

```
nutrition_sort = nutrition.sort_values(['saturated_fat_per_gram'], ascending=False)
```

```
#Display data
nutrition_sort.head(5)
```

| | | food item | weight_in_grams | saturated_fat | cholesterol | index | saturated_fat_per_gram | cholesterol_per_gram |
|-----|-----|--------------------------|-----------------|---------------|-------------|-------|------------------------|----------------------|
| 908 | 909 | BUTTER; SALTED 1 TSP | 14.0 | 7.1 | 31 | 908 | 0.507143 | |
| 909 | 710 | BUTTER; UNSALTED 1 TSP | 14.0 | 7.1 | 31 | 909 | 0.507143 | |
| 710 | 710 | BUTTER; UNSALTED 1/2 CUP | 113.0 | 57.1 | 247 | 710 | 0.505310 | |
| 913 | 709 | BUTTER; SALTED 1/2 CUP | 113.0 | 57.1 | 247 | 709 | 0.505310 | |
| 709 | 710 | BUTTER; UNSALTED 1 PAT | 5.0 | 2.5 | 11 | 913 | 0.500000 | |

Derive a new variable, cholesterol_per_gram

```
nutrition['cholesterol_per_gram'] = nutrition['cholesterol'].div(nutrition['weight_in_grams'].values)
```

```
#Display data
nutrition_chol_gram
```

| | | food item | weight_in_grams | saturated_fat | cholesterol | index | saturated_fat_per_gram | cholesterol_per_gram |
|-----|-----|-------------------------------------|-----------------|---------------|-------------|-------|------------------------|----------------------|
| 119 | 45 | EGGS; RAW; YOLK 1 YOLK | 17.0 | 1.6 | 213 | 119 | 0.094118 | 12.529412 |
| 45 | 58 | CHICKEN LIVER; COOKED 1 LIVER | 20.0 | 0.4 | 126 | 58 | 0.020000 | 6.300000 |
| 58 | 167 | BEEF LIVER; FRIED 3 OZ | 85.0 | 2.5 | 410 | 45 | 0.029412 | 4.823529 |
| 167 | 184 | EGGS; COOKED; FRIED 1 EGG | 46.0 | 1.9 | 211 | 167 | 0.041304 | 4.586957 |
| 184 | 185 | EGGS; COOKED; POACHED 1 EGG | 50.0 | 1.5 | 212 | 185 | 0.030000 | 4.240000 |
| 185 | 186 | EGGS; COOKED; HARD-COOKED 1 EGG | 50.0 | 1.6 | 213 | 186 | 0.032000 | 4.260000 |
| 186 | 189 | EGGS; COOKED; SCRAMBLED;OMLETT1 EGG | 61.0 | 2.2 | 215 | 189 | 0.036066 | 3.524590 |

961 rows × 8 columns

1. Standardize the field saturated_fat_per_gram

```
nutrition['saturated_fat_per_gram_z'] = stats.zscore(nutrition['saturated_fat_per_gram'])
```

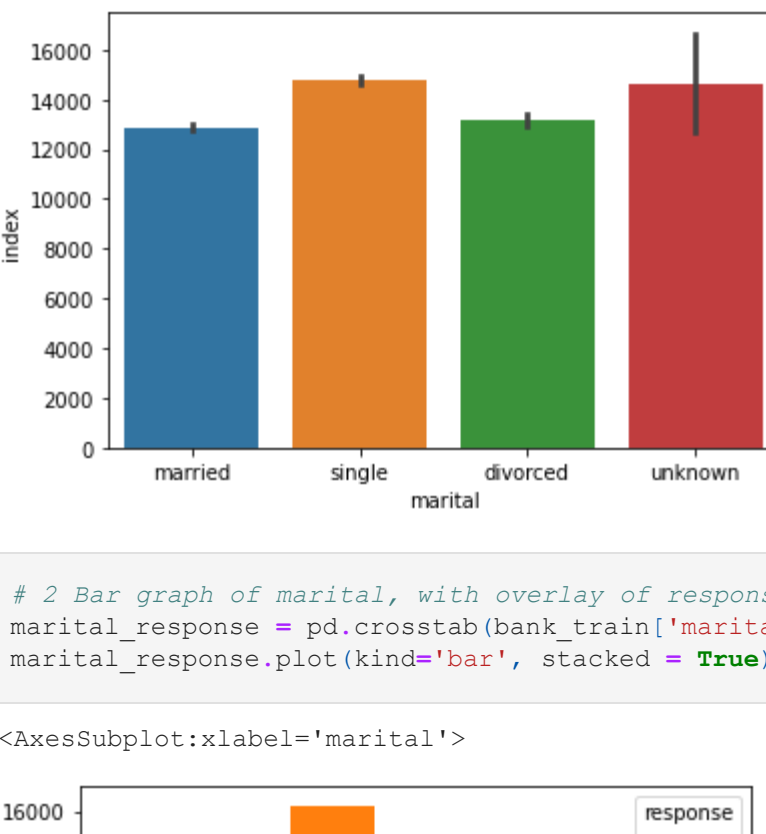
```
#Display data
nutrition
```

| | | | | | |
|-----|--|--------|------|----|-----|
| 4 | PARMESAN CHEESE GRATED 1 CUP | 100.00 | 19.1 | 79 | 4 |
| 328 | MILK CHOCOLATE CANDY; PLAIN 1 OZ | 28.35 | 5.4 | 6 | 328 |
| 74 | MAUIENSTER CHEESE 1 OZ | 28.35 | 5.4 | 27 | 74 |
| 3 | PARMESAN CHEESE; GRATED 1 OZ | 28.35 | 5.4 | 22 | 3 |
| 655 | IMITATION WHIPP TOPPING PRESERVO 1 CUP | 70.00 | 13.2 | 0 | 655 |
| 102 | BLUE CHEESE 1 OZ | 28.35 | 5.3 | 21 | 102 |
| 905 | SOYBEAN COTTONSEED OIL; HYDRONT CUP | 218.00 | 39.2 | 0 | 905 |
| 896 | SOYBEAN- COTTONSEED OIL; HYDRONT TBSP | 14.00 | 2.5 | 0 | 896 |
| 36 | SWISS CHEESE 1 OZ | 28.35 | 5.0 | 26 | 36 |
| 18 | PORK; CURED; BACON; REGULCKEDS SLICE | 19.00 | 3.3 | 16 | 18 |
| 893 | PEANUT OIL 1 TBSP | 14.00 | 2.4 | 0 | 893 |

1. standardize the feed cholesterol_per_gram

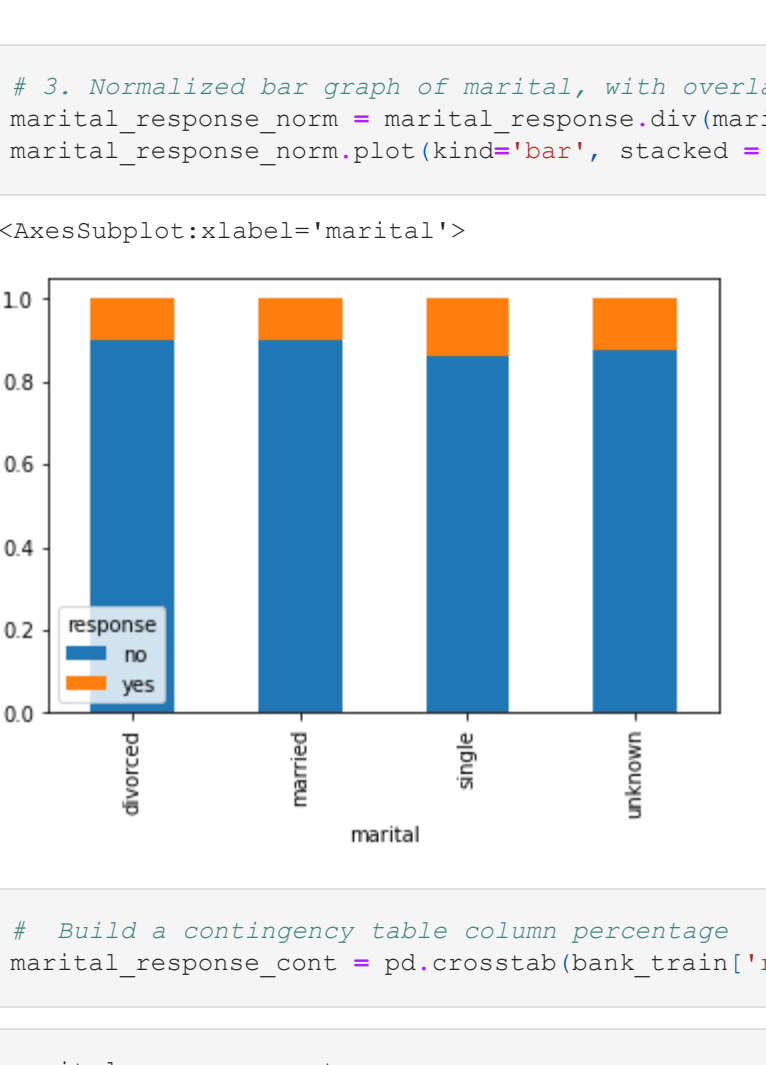
[73]:

Standardize[feeding_obs].mean.ols.



```
# 2 Bar graph of marital, with overlay of response.
marital_response = pd.crosstab(bank_train['marital'], bank_train['response'])
marital_response.plot(kind='bar', stacked = True)
```

```
<AxesSubplot: xlabel='marital'>
```



```
# 3. Normalized bar graph of marital, with overlay of response.
marital_response_norm = marital_response.div(marital_response.sum(1), axis = 0)
marital_response_norm.plot(kind='bar', stacked = True)
```

```
<AxesSubplot: xlabel='marital'>
```



```
# Build a contingency table column percentage
marital_response_cont = pd.crosstab(bank_train['marital'], bank_train['response'])
```

```
marital_response_cont
```

| marital | | | | |
|----------|------|-------|------|----|
| response | | | | |
| no | 2743 | 14579 | 6514 | 50 |
| yes | 312 | 1608 | 1061 | 7 |

```
round(marital_response_cont.div(marital_response_cont.sum(0), axis=1)*100, 1)
```

| marital | | | | |
|----------|------|------|------|------|
| response | | | | |
| no | 89.8 | 90.1 | 86.0 | 87.7 |
| yes | 10.2 | 9.9 | 14.0 | 12.3 |

```
marital_response_cont_row = pd.crosstab(bank_train['marital'], bank_train['response'])
```

```
marital_response_cont_row
```

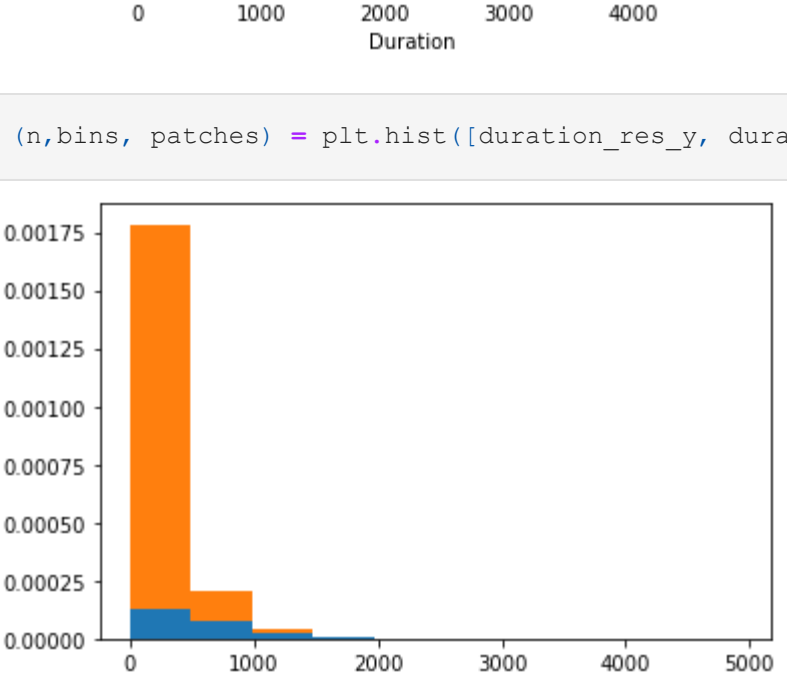
| response | | |
|----------|-------|------|
| marital | | |
| divorced | 2743 | 312 |
| married | 14579 | 1608 |
| single | 6514 | 1061 |
| unknown | 50 | 7 |

```
round(marital_response_cont_row.div(marital_response_cont_row.sum(0), axis=1)*100, 1)
```

| response | | |
|----------|------|------|
| marital | | |
| divorced | 11.5 | 10.4 |
| married | 61.0 | 53.8 |
| single | 27.3 | 35.5 |
| unknown | 0.2 | 0.2 |

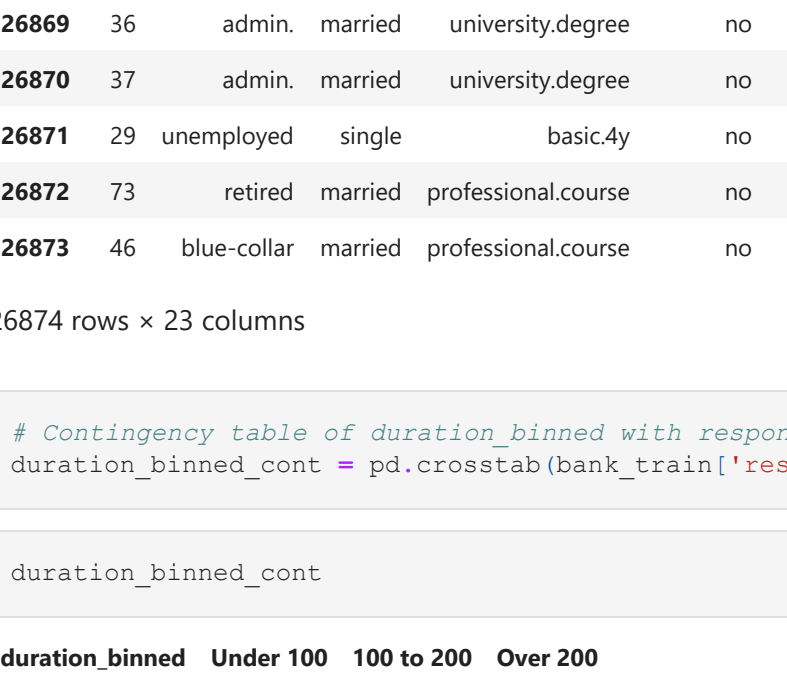
Produce the following graphs. What are the strengths and weaknesses of each graph?

```
#Histogram Plot
duration = bank_train['duration']
plt.hist((duration), rwidth=0.6)
plt.title('Histogram of Duration')
plt.xlabel('Duration')
plt.ylabel('Frequency')
plt.show()
```

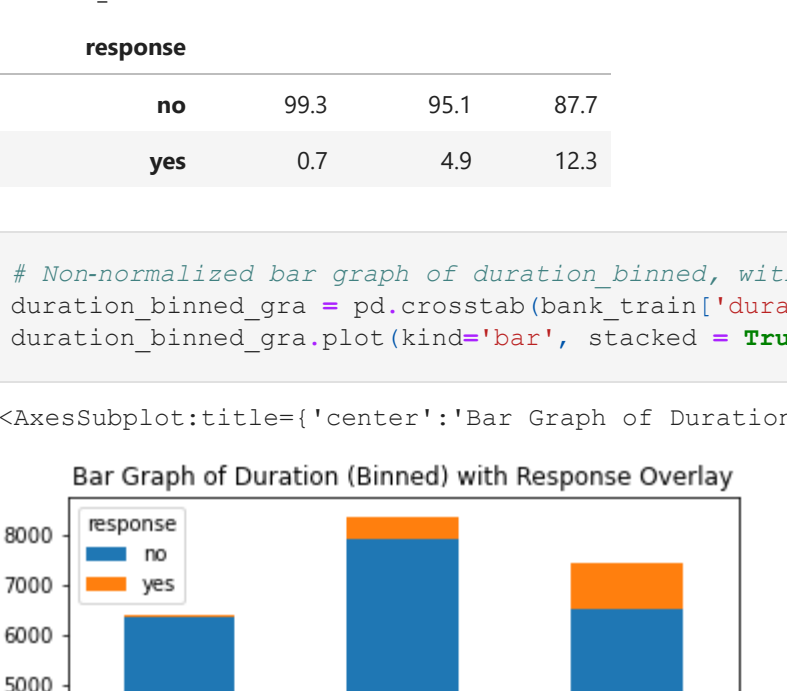


```
#Histogram of duration, with overlay of response
duration_res_y = bank_train[bank_train.response == 'yes']['duration']
duration_res_n = bank_train[bank_train.response == 'no']['duration']
```

```
plt.hist((duration_res_y, duration_res_n), bins = 10, stacked = True)
plt.legend(['Response = Yes', 'Response = No'])
plt.title('Histogram of duration with Response Overlay')
plt.xlabel('Duration')
plt.ylabel('Frequency')
plt.show()
```



```
#Normalized histogram of duration, with overlay of response.
(nbins, patches) = plt.hist((duration_res_y, duration_res_n), bins=10, stacked=True, density=True)
```



```
n_table = np.column_stack((n[0], n[1]))
```

```
n_norm = n_table / n_table.sum(axis=1)[:, None]
```

```
ourbins = np.column_stack((bins[0:10], bins[1:11]))
```

```
p1 = plt.bar(x = ourbins[:,0], height = n_norm[:,0], width = ourbins[:,1] - ourbins[:,0])
p2 = plt.bar(x = ourbins[:,0], height = n_norm[:,1], width = ourbins[:,1] - ourbins[:,0], bottom = n_norm[:,0])
plt.legend(['Response = Yes', 'Response = No'])
plt.title('Non-Normalized Histogram of Duration with Response Overlay')
plt.xlabel('Duration')
plt.ylabel('Proportion')
plt.show()
```



```
(nbins, patches) = plt.hist((duration_res_y, duration_res_n), bins=10, stacked=True, density=True)
```



```
#creating a cut of point duration_binned
bank_train['duration_binned'] = pd.cut(x=bank_train['duration'], bins=[0,100,200,400], labels=['Under 100', '100 to 200', 'Over 200'])
```

```
bank_train
```

| | age | job | marital | education | default | housing | loan | contact | month | day_of_week | ... | previous | outcome | |
|-------|-----|-------------|---------|---------------------|---------|---------|------|-----------|-----------|-------------|-----|----------|-------------|-------------|
| 0 | 56 | services | married | highschool | unknown | no | no | telephone | may | mon | .. | 0 | nonexistent | |
| 1 | 41 | blue-collar | married | unknown | unknown | no | no | telephone | may | mon | .. | 0 | nonexistent | |
| 2 | 41 | blue-collar | married | unknown | unknown | no | yes | no | telephone | may | mon | .. | 0 | nonexistent |
| 3 | 25 | services | single | highschool | no | yes | no | telephone | may | mon | .. | 0 | nonexistent | |
| 4 | 29 | blue-collar | single | highschool | no | no | yes | telephone | may | mon | .. | 0 | nonexistent | |
| ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | |
| 26869 | 36 | admin. | married | university.degree | no | no | no | cellular | nov | fri | .. | 0 | nonexistent | |
| 26870 | 37 | admin. | married | university.degree | no | yes | no | cellular | nov | fri | .. | 0 | nonexistent | |
| 26871 | 29 | unemployed | single | basic-4y | no | yes | no | cellular | nov | fri | .. | 1 | success | |
| 26872 | 73 | retired | married | professional.course | no | yes | no | cellular | nov | fri | .. | 0 | nonexistent | |
| 26873 | 46 | blue-collar | married | professional.course | no | no | no | cellular | nov | fri | .. | 0 | nonexistent | |

26874 rows x 23 columns

```
# Contingency table of duration_binned with response
duration_binned_cont = pd.crosstab(bank_train['response'], bank_train['duration_binned'])
```

```
duration_binned_cont
```

| duration_binned | | | |
|-----------------|------|------|------|
| response | | | |
| no | 6360 | 7945 | 6517 |
| yes | 46 | 406 | 913 |

```
# count of column percentage
round(duration_binned_cont.div(duration_binned_cont.sum(0), axis=1)*100, 1)
```

| duration_binned | | | |
|-----------------|------|------|------|
| response | | | |
| no | 99.3 | 95.1 | 87.7 |
| yes | 0.7 | 4.9 | 12.3 |

```
# Non-normalized bar graph of duration_binned, with response overlay.
duration_binned_gra = pd.crosstab(bank_train['duration_binned'], bank_train['response'])
duration_binned_gra.plot(kind='bar', stacked = True, title = "Bar Graph of Duration (Binned) with Response Overlay")
```

```
<AxesSubplot: title='(center): Bar Graph of Duration (Binned) with Response Overlay', xlabel='duration_binned'>
```

