# Sipna College of Engineering & Technology, Amravati.
## Department of Computer Science & Engineering
### Session 2022-2023

**Branch :- Computer Sci. & Engg.**                    **Class :- Final Year**
**Subject :-Artificial Intelligence and Machine Learning**    **Sem  :- VIII**
**Teacher Manual**

| PRACTICAL NO 3 |
|---|

**AIM**:   To Understand and implement the concept of un-supervised learning

**S/W REQUIRED:** Python
**DATA SET USED:** iris.csv

## Unsupervised Learning:-

Unsupervised learning, also known as unsupervised machine learning, uses machine learning algorithms to analyze and cluster unlabeled datasets

Unsupervised learning models are utilized for three main tasks—clustering, association, and dimensionality reduction. Below we'll define each learning method and highlight common algorithms and approaches to conduct them effectively.

Clustering is a data mining technique which groups unlabelled data based on their similarities or differences. Clustering algorithms are used to process raw, unclassified data objects into groups represented by structures or patterns in the information. Clustering algorithms can be categorized into a few types, specifically exclusive, overlapping, hierarchical, and probabilistic.

**Implementation:**
```python
# Importing the Libraries
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns

%matplotlib inline
```

```python
# Importing the Dataset
try:
    data = pd.read_csv("../input/Wholesale customers data.csv")
    data.drop(labels=(['Channel','Region']),axis=1,inplace=True)
    print('Wholesale customers has {} samples with {} features each'.format(*data.shape))
except:
    print('Sorry! Dataset could not be loaded.')
```

**O/P:-** Wholesale customers has 440 samples with 6 features each

```python
data.head()

# Display a brief description of the overall dataset
data.describe()
```

|  | Fresh | Milk | Grocery | Frozen | Detergents_Paper | Delicassen |
|---|---|---|---|---|---|---|
| count | 440.000000 | 440.000000 | 440.000000 | 440.000000 | 440.000000 | 440.000000 |
| mean | 12000.297727 | 5796.265909 | 7951.277273 | 3071.931818 | 2881.493182 | 1524.870455 |
| std | 12647.328865 | 7380.377175 | 9503.162829 | 4854.673333 | 4767.854448 | 2820.105937 |
| min | 3.000000 | 55.000000 | 3.000000 | 25.000000 | 3.000000 | 3.000000 |
| 25% | 3127.750000 | 1533.000000 | 2153.000000 | 742.250000 | 256.750000 | 408.250000 |
| 50% | 8504.000000 | 3627.000000 | 4755.500000 | 1526.000000 | 816.500000 | 965.500000 |
| 75% | 16933.750000 | 7190.250000 | 10655.750000 | 3554.250000 | 3922.000000 | 1820.250000 |
| max | 112151.000000 | 73498.000000 | 92780.000000 | 60869.000000 | 40827.000000 | 47943.000000 |

*# Display complete information of the data frame*
data.info()

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 440 entries, 0 to 439
Data columns (total 6 columns):
Fresh               440 non-null int64
Milk                440 non-null int64
Grocery             440 non-null int64
Frozen              440 non-null int64
Detergents_Paper    440 non-null int64
Delicassen          440 non-null int64
dtypes: int64(6)
memory usage: 20.7 KB
```

*# Select three indices of your choice you wish to sample from the dataset*
indices = [22,154,398]

*# Create a DataFrame of the chosen samples*
samples = pd.DataFrame(data.loc[indices], columns=data.keys()).reset_index(drop=True)
print("Chosen samples of wholesale customers dataset:")
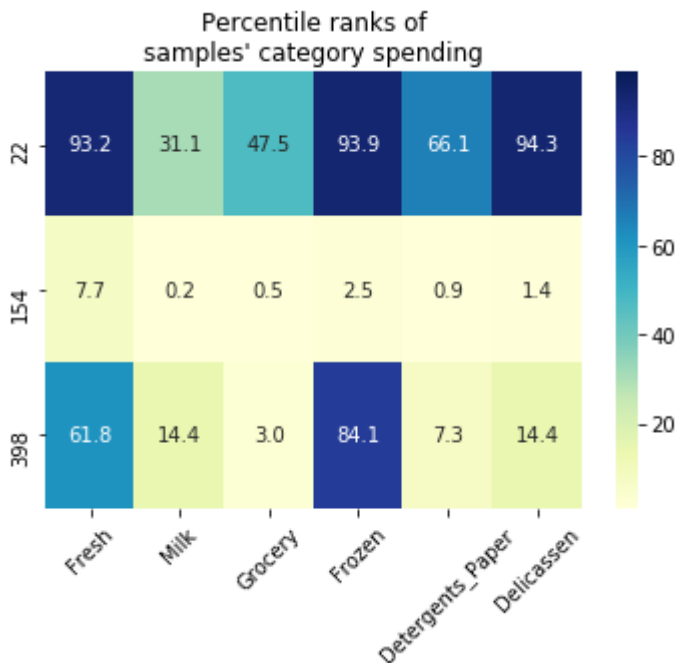display(samples)

```
Chosen samples of wholesale customers dataset:
```

|  | Fresh | Milk | Grocery | Frozen | Detergents_Paper | Delicassen |
|---|---|---|---|---|---|---|
| 0 | 31276 | 1917 | 4469 | 9408 | 2381 | 4334 |
| 1 | 622 | 55 | 137 | 75 | 7 | 8 |
| 2 | 11442 | 1032 | 582 | 5390 | 74 | 247 |

[Type text]

```python
# look at percentile ranks
#pcts = 100. * data.rank(axis=0, pct=True).iloc[indices].round(decimals=3)
pcts = 100. * data.rank(axis=0, pct=True).iloc[indices].round(decimals=3)
# visualize percentiles with heatmap

sns.heatmap(pcts, annot=True, vmin=1, vmax=99, fmt='.1f', cmap='YlGnBu')
plt.title('Percentile ranks of\nsamples\' category spending')
plt.xticks(rotation=45, ha='center');
```

Percentile ranks of
samples' category spending

|   | Fresh | Milk | Grocery | Frozen | Detergents_Paper | Delicassen |
|---|-------|------|---------|--------|------------------|------------|
| 22 | 93.2 | 31.1 | 47.5 | 93.9 | 66.1 | 94.3 |
| 154 | 7.7 | 0.2 | 0.5 | 2.5 | 0.9 | 1.4 |
| 398 | 61.8 | 14.4 | 3.0 | 84.1 | 7.3 | 14.4 |

```python
# Import libraries for Decision Tree Regressor
from sklearn.tree import DecisionTreeRegressor
from sklearn.model_selection import train_test_split

# Remove column Milk
new_data = data.drop('Milk',axis=1)

# Split the data into training and testing sets(0.25) using the given feature as the target
# Set a random state.
X_train, X_test, y_train, y_test = train_test_split(new_data, data['Milk'], test_size=0.25, random_state=1)

# Create a decision tree regressor and fit it to the training set
regressor = DecisionTreeRegressor(random_state=1)
regressor.fit(X_train, y_train)

# Report the score of the prediction using the testing set
score = regressor.score(X_test, y_test)
print(score)
```

**Accuracy:-** 0.515849943807

**CONCLUSION:** Thus we have implemented the concept of un-supervised learning.

[Type text]

(**As you can see, we attempted to predict Milk using the other features in the dataset and the score ended up being 0.515. At this initial stage we might say that this feature is somewhat difficult to predict because the score is around the halfway point of possible scores. Remember that R^2 goes from 0 to 1. This might indicate that it could be an important feature to consider.**)