

# K-Prototype and K-Modes

## K-Prototype

Import Library

1. Lakukan langkah yang sama seperti minggu sebelumnya untuk meng-import library Numpy dan Pandas, numpy, plotnine, module for k-prototype cluster.

```
In [ ]: # Import module for data manipulation
import pandas as pd

# Import module for linear algebra
import numpy as np

import re

# Import module for data visualization
from plotnine import *
import plotnine

# Import module for k-prototype cluster
from knodes.kprototypes import KPrototypes

# Import matplotlib
from matplotlib.figure import Figure
import matplotlib.pyplot as plt

# Ignore warnings
warnings.filterwarnings('ignore', category=FutureWarning)

C:\Users\Darren\anaconda3\lib\site-packages\scipy\_init_.py:146: UserWarning: A NumPy version >=1.16.5 and <1.23.0 is required for this version of SciPy (detected version 1.24.3)

In [ ]: # Format scientific notation from Pandas
pd.set_option("display.float_format", lambda x: '%.3f' % x)

Import Data
```

2. Gunakan dataset 10000 Sales Records.csv dan masukkan ke dalam dataframe menggunakan Pandas. Berikan nama df (Petunjuk: gunakan delimiter=).

3. Tampilkan informasi dataset dan isi data.

```
In [ ]: df = pd.read_csv(r'D:\SEMESTER 4\IS411 Data Modelling\LAB\IS411_C-HY_00000054804_Christopher_Darren_Week-12\10000 Sales Records.csv', delimiter=',')

In [ ]: print("Dimension data: {} rows and {} columns".format(len(df), len(df.columns)))
df.head()

Dimension data: 10000 rows and 14 columns
```

	Region	Country	Item Type	Sales Channel	Order Priority	Order Date	Order ID	Ship Date	Units Sold	Unit Price	Unit Cost	Total Revenue	Total Cost	Total Profit
0	Sub-Saharan Africa	Chad	Office Supplies	Online	L	1/27/2011	292494523	2/12/2011	4484	651.210	524.960	2920025.640	2353920.640	566105.000
1	Europe	Latvia	Beverages	Online	C	12/28/2015	361829549	1/29/2016	1075	47.450	31.790	51008.750	34174.250	16834.500
2	Middle East and North Africa	Pakistan	Vegetables	Offline	C	1/13/2011	141515767	2/1/2011	6515	154.000	90.930	1003700.900	592408.950	411291.950
3	Sub-Saharan Africa	Democratic Republic of the Congo	Household	Online	C	9/11/2012	500364005	10/6/2012	7683	668.270	502.540	5134318.410	3861014.820	1273303.590
4	Europe	Czech Republic	Beverages	Online	C	10/27/2015	127481591	12/5/2015	3491	47.450	31.790	165647.950	110978.890	54669.060

4. Inspect the data type

```
In [ ]: df.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 10000 entries, 0 to 9999
Data columns (total 14 columns):
 #   Column                Non-Null Count  Dtype
---  --
 0   Region                10000 non-null object
 1   Country               10000 non-null object
 2   Item Type             10000 non-null object
 3   Sales Channel         10000 non-null object
 4   Order Priority        10000 non-null object
 5   Order Date           10000 non-null object
 6   Order ID             10000 non-null int64
 7   Ship Date            10000 non-null object
 8   Units Sold           10000 non-null int64
 9   Unit Price           10000 non-null float64
10   Unit Cost            10000 non-null float64
11   Total Revenue        10000 non-null float64
12   Total Cost           10000 non-null float64
13   Total Profit         10000 non-null float64
dtypes: float64(5), int64(2), object(7)
memory usage: 1.1+ MB

5. Inspect the categorical variables

In [ ]: df.select_dtypes('object').nunique()

Out [ ]: Region          7
Country        165
Item Type       12
Sales Channel   2
Order Priority   4
Order Date     2691
Ship Date      2719
dtype: int64

6. Inspect the numerical variables.

In [ ]: df.describe()

Out [ ]:
```

	Order ID	Units Sold	Unit Price	Unit Cost	Total Revenue	Total Cost	Total Profit
count	10000.000	10000.000	10000.000	10000.000	10000.000	10000.000	10000.000
mean	549871874.366	5002.856	268.143	198.807	133395.131	93805.784	39509.347
std	260787951.133	2873.246	217.944	176.446	1466026.174	1145914.069	377554.961
min	100380161.000	2.000	9.330	6.920	167.940	124.560	43.380
25%	321806669.000	2530.750	109.280	56.670	288561.078	164785.530	98329.140
50%	548665305.000	4562.000	205.700	117.110	800561.210	481605.840	288999.020
75%	799598103.500	7472.000	437.200	364.960	1819143.390	1183021.520	566422.708
max	999994232.000	10000.000	668.270	524.960	6680028.820	5241725.600	1788178.380

7. Check missing value.

```
In [ ]: df.isna().sum()

Out [ ]: Region          0
Country          0
Item Type        0
Sales Channel    0
Order Priority    0
Order Date      0
Ship Date       0
Units Sold      0
Unit Price      0
Unit Cost       0
Total Revenue   0
Total Cost      0
Total Profit    0
dtype: int64

Explanatory Data Analysis
```

8. Distribution of region

```
In [ ]: df_region = pd.DataFrame(df['Region'].value_counts()).reset_index()
df_region['Percentage'] = df_region['Region'] / df_region['value_counts().sum()
df_region.rename(columns={'index': 'Region', 'Region': 'Total', 'inplace=True'})
df_region = df_region.sort_values('Total', ascending=True).reset_index(drop=True)
df_region

Out [ ]:
```

	Region	Total	Percentage
0	North America	215	0.021
1	Australia and Oceania	797	0.080
2	Central America and the Caribbean	1019	0.102
3	Middle East and North Africa	1264	0.126
4	Asia	1469	0.147
5	Sub-Saharan Africa	2603	0.260
6	Europe	2633	0.263

```
In [ ]: # The DataFrame
df_region = df_region.groupby('Region').agg({
    'Region': 'count',
    'Units Sold': 'mean',
    'Total Revenue': 'mean',
    'Total Cost': 'mean',
    'Total Profit': 'mean'
})
df_region.rename(columns={'Region': 'Total'}).reset_index().sort_values('Total', ascending=True)

In [ ]: df_region

Out [ ]:
```

	Region	Total	Units Sold	Total Revenue	Total Cost	Total Profit
5	North America	215	5373.358	1599778.805	1097008.967	462769.838
1	Australia and Oceania	797	4966.769	1317192.267	910578.451	406613.816
2	Central America and the Caribbean	1019	5061.063	1369509.041	973672.093	395836.948
4	Middle East and North Africa	1264	5116.219	1357304.882	953884.190	403420.693
0	Asia	1469	5015.488	1365082.080	965115.963	399966.097
6	Sub-Saharan Africa	2603	4967.808	1287190.079	933155.469	384034.611
3	Europe	2633	4920.385	1222207.396	932138.169	390049.226

9. Data Visualization

```
In [ ]: plotnine.options.figure_size = (8, 4.8)

(
    ggplot(data=df_region) +
    geom_bar(aes(x='Region',
                 y='Total', fill=np.where(df_region['Region'] == 'Asia', '#981220', '#00797c')),
             stat='identity') +
    geom_text(aes(x='Region',
                  y='Total',
                  label='Total'),
              size=16,
              nudge_y=120) +
    labs(title='Region that has the highest purchases') +
    xlab('Region') +
    ylab('Frequency') +
    scale_x_discrete(limits = df_region['Region'].tolist()) +
    theme_minimal()
)

Region that has the highest purchases
```

10. Distribution of item type

```
In [ ]: # Order the index of cross tabulation
order_region = df_region['Region'].tolist()
order_region.append('All')
order_region

Out [ ]: ['North America',
'Australia and Oceania',
'Central America and the Caribbean',
'Middle East and North Africa',
'Asia',
'Sub-Saharan Africa',
'Europe',
'All']

In [ ]: df_item = pd.crosstab([df['Region'], df['Item Type']], margins=True).reindex(order_region, axis=0).reset_index()
df_item.index.name = None
df_item.columns.name = None
df_item

Out [ ]:
```

	Region	Baby Food	Beverages	Cereal	Clothes	Cosmetics	Fruits	Household	Meat	Office Supplies	Personal Care	Snacks	Vegetables	All
0	North America	21	20	16	21	20	15	20	17	20	17	16	12	215
1	Australia and Oceania	65	50	69	77	75	55	78	61	50	76	72	69	797
2	Central America and the Caribbean	74	92	77	84	77	81	104	75	94	82	89	90	1019
3	Middle East and North Africa	105	96	104	111	99	104	128	101	103	112	95	106	1264
4	Asia	132	108	121	116	125	111	116	114	132	137	120	137	1469
5	Sub-Saharan Africa	235	220	211	229	203	230	218	207	207	223	221	199	2603
6	Europe	210	196	227	234	235	199	211	223	231	241	203	233	2633
7	All	842	782	825	874	834	795	875	798	837	888	816	836	10000

Data Pre-processing

11. Data pre-processing

```
In [ ]: df.drop(['Country', 'Order Date', 'Order ID', 'Ship Date'], axis=1, inplace=True)

In [ ]: print("Dimension data: {} rows and {} columns".format(len(df), len(df.columns)))
df.head()

Dimension data: 10000 rows and 10 columns
```

	Region	Item Type	Sales Channel	Order Priority	Units Sold	Unit Price	Unit Cost	Total Revenue	Total Cost	Total Profit
0	Sub-Saharan Africa	Office Supplies	Online	L	4484	651.210	524.960	2920025.640	2353920.640	566105.000
1	Europe	Beverages	Online	C	1075	47.450	31.790	51008.750	34174.250	16834.500
2	Middle East and North Africa	Vegetables	Offline	C	6515	154.000	90.930	1003700.900	592408.950	411291.950
3	Sub-Saharan Africa	Household	Online	C	7683	668.270	502.540	5134318.410	3861014.820	1273303.590
4	Europe	Beverages	Online	C	3491	47.450	31.790	165647.950	110978.890	54669.060

Cluster Analysis

12. Get the position of categorical columns

```
In [ ]: catColumnsPos = [df.columns.get_loc(col) for col in list(df.select_dtypes('object').columns)]
print('Categorical columns position : {}'.format(catColumnsPos))

Categorical columns : ['Region', 'Item Type', 'Sales Channel', 'Order Priority']
Categorical columns position : [0, 1, 2, 3]

13. Convert dataframe to matrix
```

```
In [ ]: dfMatrix = df.to_numpy()

In [ ]: dfMatrix

array([[ 'Sub-Saharan Africa', 'Office Supplies', 'Online', ...,
        2920025.64, 2353920.64, 566105.0],
       [ 'Europe', 'Beverages', 'Online', ...,
        51008.75, 34174.25,
        16834.5],
       [ 'Middle East and North Africa', 'Vegetables', 'Offline', ...,
        1003700.9, 592408.95, 411291.95],
       [ 'Sub-Saharan Africa', 'Vegetables', 'Offline', ...,
        3861014.82, 2353920.64, 1273303.59],
       [ 'Sub-Saharan Africa', 'Household', 'Online', ...,
        5134318.41, 3861014.82, 1273303.59],
       [ 'Europe', 'Beverages', 'Online', ...,
        165647.95, 110978.89, 54669.06],
       [ 'Asia', 'Snacks', 'Offline', ...,
        595881.38, 35175.84, 19905.54]],
      dtype=object)
```

14. Choose optima K using Elbow method

```
In [ ]: # Choosing optimal K
cost = []
for cluster in range(1, 10):
    try:
        kprototype = KPrototypes(n_jobs=-1, n_clusters=cluster, init='Huang', random_state=0)
        kprototype.fit_predict(dfMatrix, categorical=catColumnsPos)
        cost.append(kprototype.cost_)
        print('Cluster initialization: {}'.format(cluster))
    except:
        break

Cluster initialization: 1
Cluster initialization: 2
Cluster initialization: 3
Cluster initialization: 4
Cluster initialization: 5

In [ ]: # Converting the result into a dataframe and plotting them
df_cost = pd.DataFrame({'Cluster': range(1,6), 'Cost': cost})
df_cost.head()

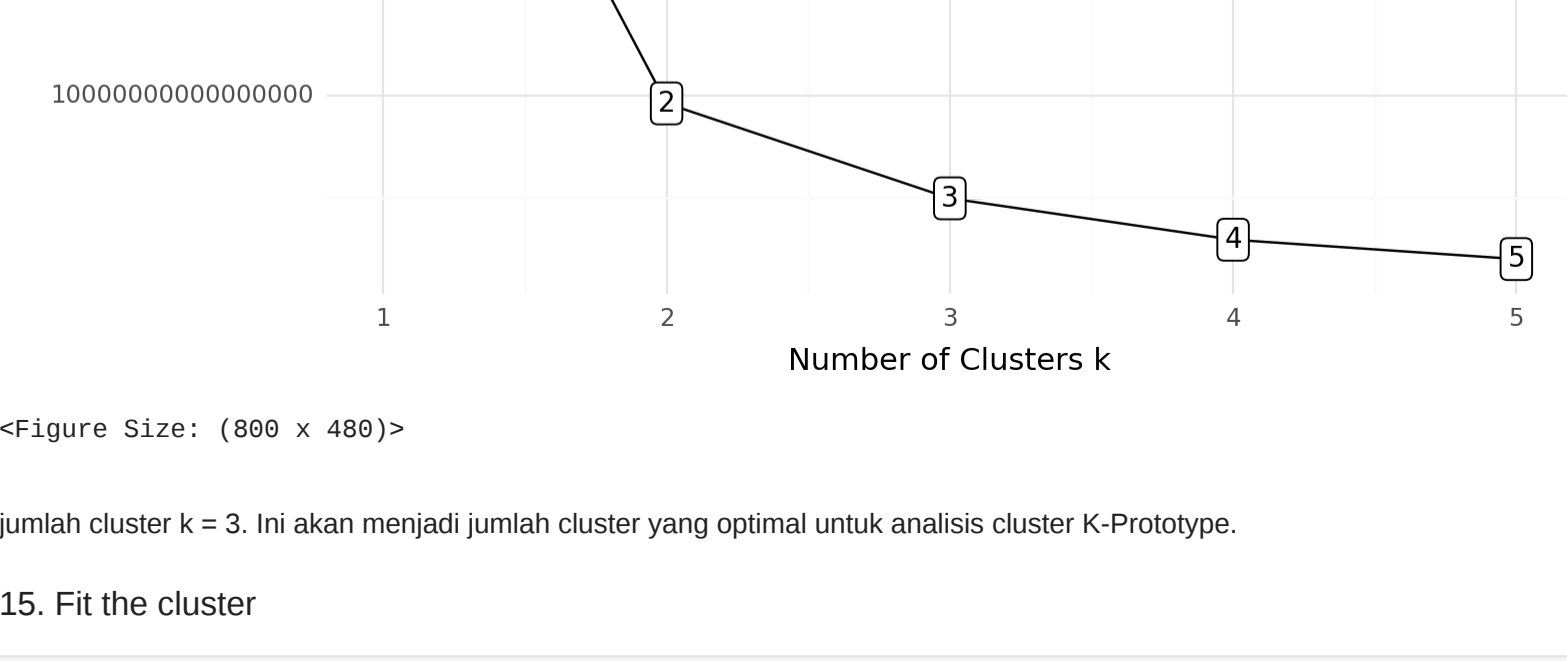
Out [ ]:
```

	Cluster	Cost
0	1	36016173492135216.000
1	2	9627968062172950.000
2	3	4960707512782198.000
3	4	2927457054391281.500
4	5	1975342819318915.500

```
In [ ]: # Data Visualization
plotnine.options.figure_size = (8, 4.8)

(
    ggplot(data=df_cost) +
    geom_line(aes(x='Cluster',
                  y='Cost')) +
    geom_point(aes(x='Cluster',
                   y='Cost')) +
    geom_label(aes(x='Cluster',
                   y='Cost',
                   label='Cluster'),
              size=10,
              nudge_y=1000) +
    labs(title='Optimal number of cluster with Elbow Method') +
    xlab('Number of clusters K') +
    ylab('Cost') +
    theme_minimal()
)
```

Optimal number of cluster with Elbow Method



jumlah cluster k = 3. Ini akan menjadi jumlah cluster yang optimal untuk analisis cluster K-Prototype.

15. Fit the cluster

```
In [ ]: # Fit the cluster
kprototype = KPrototypes(n_jobs=-1, n_clusters=3, init='Huang', random_state=0)
kprototype.fit_predict(dfMatrix, categorical=catColumnsPos)

array([2, 1, 1, ..., 1, 0, 1], dtype=uint8)
```

16. Cluster centroid

```
In [ ]: kprototype.cluster_centroids_

array([[ 'Sub-Saharan Africa', 'Office Supplies', 'Online', ...,
        593.5265126998487, '457.7854957983237',
        4632769.54755464, '1385921.134847464', '1983839.424867272',
        'Europe', 'Household', 'Offline', ...,
        105.7960808293876, '163.29510972380268', '105.7960808293876',
        467709.4517983619, '281142.954697418', '186566.49795961852',
        'Sub-Saharan Africa', 'Personal Care', 'Online', 'C',
        16983.278421846355, '384.26403810793', '276.69765819222703',
        1995888.1232729563, '1388539.5273834165', '615348.5959695351',
        'Europe', 'Cosmetics', 'Online', 'H', ...,
        1995888.1232729563, '1388539.5273834165', '615348.5959695351'],
      dtype='<U32')

In [ ]: # Check the iteration of the clusters created
kprototype.n_iter_

Out [ ]: 7
```

18. Check the cost of the clusters created

```
In [ ]: # Check the cost of the clusters created
kprototype.cost_

Out [ ]: 4960707512782198.0

19. Add the cluster to the dataframe
```

```
In [ ]: # Add the cluster to the dataframe
df['cluster_id'] = kprototype.labels_

df_region = pd.DataFrame(df['Region'].value_counts()).reset_index()
df_region['Percentage'] = df_region['Region'] / df_region['value_counts().sum()
df_region.rename(columns={'index': 'Region', 'Region': 'Total', 'inplace=True'})
df_region = df_region.sort_values('Total', ascending=True).reset_index(drop=True)
df_region

Out [ ]:
```

	Region	Total	Percentage
0	North America	215	0.021
1	Australia and Oceania	797	0.080
2	Central America and the Caribbean	1019	0.102
3	Middle East and North Africa	1264	0.126
4	Asia	1469	0.147
5	Sub-Saharan Africa	2603	0.260
6	Europe	2633	0.263

```
In [ ]: # Add the cluster to the dataframe
df['Cluster Labels'] = kprototype.labels_
df['Segment'] = df['Cluster Labels'].map({0: 'First', 1: 'Second', 2: 'Third'})

# Order the cluster
df['Segment'] = df['Segment'].astype('category')
df['Segment'] = df['Segment'].cat.reorder_categories(['First', 'Second', 'Third'])
df.head()

Out [ ]:
```

	Region	Item Type	Sales Channel	Order Priority	Units Sold	Unit Price	Unit Cost	Total Revenue	Total Cost	Total Profit	cluster_id	Cluster Labels	Segment
0	Sub-Saharan Africa	Office Supplies	Online	L	4484	651.210	524.960	2920025.640	2353920.640	566105.000	2	2	Third
1	Europe	Beverages	Online	C	1075	47.450	31.790	51008.750	34174.250	16834.500	1	1	Second
2	Middle East and North Africa	Vegetables	Offline	C	6515	154.000	90.930	1003700.900	592408.950	411291.950	1	1	Second
3	Sub-Saharan Africa	Household	Online	C	7683	668.270	502.540	5134318.410	3861014.820	1273303.590	0	0	First
4	Europe	Beverages	Online	C	3491	47.450	31.790	165647.950	110978.890	54669.060	1	1	Second

20. Cluster interpretation

```
In [ ]: df.rename(columns={'Cluster Labels': 'Total', 'inplace=True'})
df.groupby('Segment').agg({
    'Total': 'count',
    'Region': lambda x: x.value_counts().index[0],
    'Item Type': lambda x: x.value_counts().index[0],
    'Sales Channel': lambda x: x.value_counts().index[0],
    'Order Priority': lambda x: x.value_counts().index[0],
    'Units Sold': 'mean',
    'Unit Price': 'mean',
    'Total Revenue': 'mean',
    'Total Cost': 'mean',
    'Total Profit': 'mean'
}).reset_index()

Out [ ]:
```

	Segment	Total	Region	Item Type	Sales Channel	Order Priority	Units Sold	Unit Price	Total Revenue	Total Cost	Total Profit
0	First	1190	Europe	Household	Offline	L	7504.366	524.960	4027760.548	3556121.123	1953639.424
1	Second	6381	Sub-Saharan Africa	Personal Care	Online	C	4046.669	163.295	467709.452	281142.955	186566.497
2	Third	2429	Europe	Cosmetics	Online	H	6083.275	384.265	1995888.123	1380539.527	615348.596

## K-Modes Clustering

1. Import Library

```
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt

2. gunakan dataset dummy dan lihat hasilnya

In [ ]: data = pd.read_csv(r'D:\SEMESTER 4\IS411 Data Modelling\LAB\IS411_C-HY_00000054804_Christopher_Darren_Week-12\data.csv')

Out [ ]:
```

	Individual	Skill
0	1	C
1	1	Python
2	1	Java
3	1	Haskell
4	2	Python
5	2	React
6	2	JS
7	2	PHP
8	3	C++
9	3	JS
10	3	Flutter
11	3	Android
12	3	iOS
13	4	Fortran
14	4	Pascal
15	4	.NET
16	4	C#
17	4	C
18	4	MATLAB

3. Kita akan memberikan tanda kepada individual yang memiliki skill dengan cara :

```
In [ ]: one_hot_df = data.copy()
for i, name in enumerate(data['Skill'].unique()):
    one_hot_df[name] = 0
def set_product(x):
    x[str('Skill')] = 1
    return x

one_hot_df = one_hot_df.apply(set_product, axis=1)
one_hot_df = one_hot_df.groupby(['Individual'], axis=1)
one_hot_df

Out [ ]:
```

	Individual	Python	Java	Haskell	React	JS	PHP	C++	JS	Flutter	iOS	Fortran	Pascal	.NET	C#	MATLAB
0	1	1	1	1	1	0	0	0	0	0	0	0	0	0	0	0
1	2	0	1	1	0	0	1	1	1							