



One can clearly see that K-Means provides the best results when clustering the given dataset. DBSCAN isn't represented in this plot since it performs really poorly as seen above. In the further process of this project and to perform the cluster analysis we will therefore use the **K-Means method with k=3 clusters**.

4. Cluster Analysis

4.1. Assess Feature Importance

```
In [39]: def get_ss_variables(df):  
    """Get the SS for each variable"""  
    ss_vars = df.var() * (df.count() - 1)  
    return ss_vars  
  
def r2_variables(df, labels):  
    """Get the R^2 for each variable"""  
    ss_vars = get_ss_variables(df)  
    ssw_vars = np.sum(df.groupby(labels).apply(get_ss_variables))  
    return ss_vars/ssw_vars
```

```
In [40]: feature_importance = pd.concat((df_scaled, pd.Series(km_labels, name='labels')), axis=1)  
feature_importance = r2_variables(feature_importance, 'labels').sort_values(ascending=False).drop(['labels'])  
feature_importance.to_frame().plot.bar(figsize=(15,7))  
plt.title("Feature Importance", fontsize=21)  
plt.legend().remove()  
plt.xlabel("Features", fontsize=13)  
plt.ylabel("R^2", fontsize=13)  
plt.show()
```



4.2. Overview

```
In [41]: cluster_analysis = clusters.copy()
```

```
In [42]: # replace cluster labels from 0 to 1, 1 to 2 and 2 to 3  
cluster_analysis['labels'] = cluster_analysis.labels + 1
```

```
In [43]: cluster_analysis.labels.value_counts()
```

```
Out[43]: 3    3527  
         1    3486  
         2    2987  
         Name: labels, dtype: int64
```

```
In [44]: cluster_means = cluster_analysis.groupby("labels").mean()  
cluster_means
```

	Daysvuss	Age	Edu	Income	Freq	Recency	Monetary	LTV	Perdeal	Dryed	Sweetest
labels											
1	904.106139	66.690763	16.833333	98942.249570	27.907344	50.090361	1339.666667	525.915950	5.655766	46.444636	
2	896.515902	28.707399	15.643120	40448.831938	4.100435	80.407097	89.354536	1.138601	57.272179	29.516237	
3	893.510916	45.659200	17.574142	66149.769492	10.419053	59.335696	365.345052	72.006805	37.761270	71.946697	

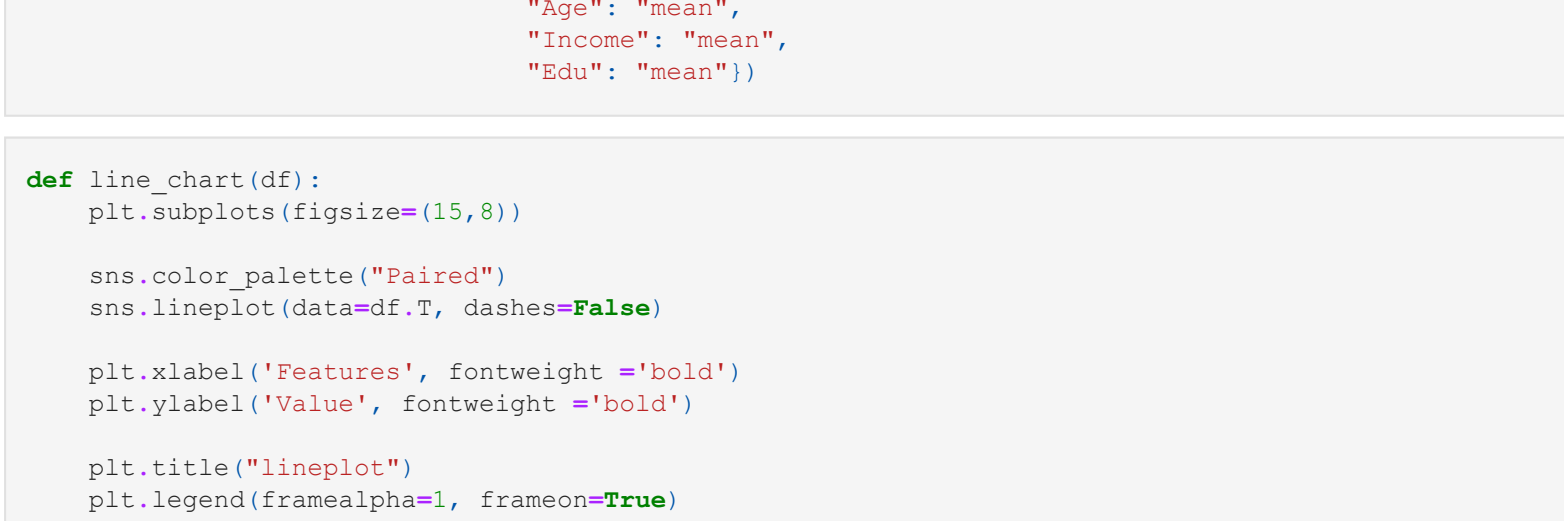
```
In [45]: cluster_means = cluster_means.append(df_before.mean(), ignore_index=True)  
cluster_means["Cluster"] = ["Cluster 1", "Cluster 2", "Cluster 3", "Average"]
```

```
In [46]: #ceil(len(features)) / 3  
fig, axes = plt.subplots(5, ceil(len(features) / 3), figsize=(21, 17))  
  
metric_features = df_before.columns.to_list()  
for ax, feat in zip(axes.flatten(), metric_features):  
    sns.barplot(x="Cluster", y=feat, data=cluster_means.sort_values(feat, ascending=False),  
               palette=palette("coolwarmblue" if feat=="Average" else "navy" for x in cluster_means.sort_values(feat, ascending=False).cluster),  
               ax=ax, set_xlabel='')  
    sns.despine(right=True)
```

```
#title = "Cluster Comparison to Mean of Feature"  
#plt.suptitle(title)
```

```
# Rotating X-axis labels  
axes.flatten()[0].tick_params(axis='x', labelrotation=90)  
#axes.flatten()[1].tick_params(axis='x', labelrotation=90)  
axes.flatten()[2].tick_params(axis='x', labelrotation=90)  
plt.subplots_adjust(wspace=0.3, hspace=0.7)
```

```
plt.show()
```



4.3. Socio-Demographic Perspective

```
In [47]: child = df['Child']  
socio_demo = cluster_analysis.copy()  
socio_demo = socio_demo.merge(child, how='outer', left_index=True, right_index=True)  
socio_demo = socio_demo.groupby("labels").agg({"labels": "count",  
        "Age": "mean",  
        "Income": "mean",  
        "Edu": "mean",  
        "WebPurchase": "mean",  
        "WebVisit": "mean",  
        "Child": "mean" # non metric feature, but appended to table for comparison  
})
```

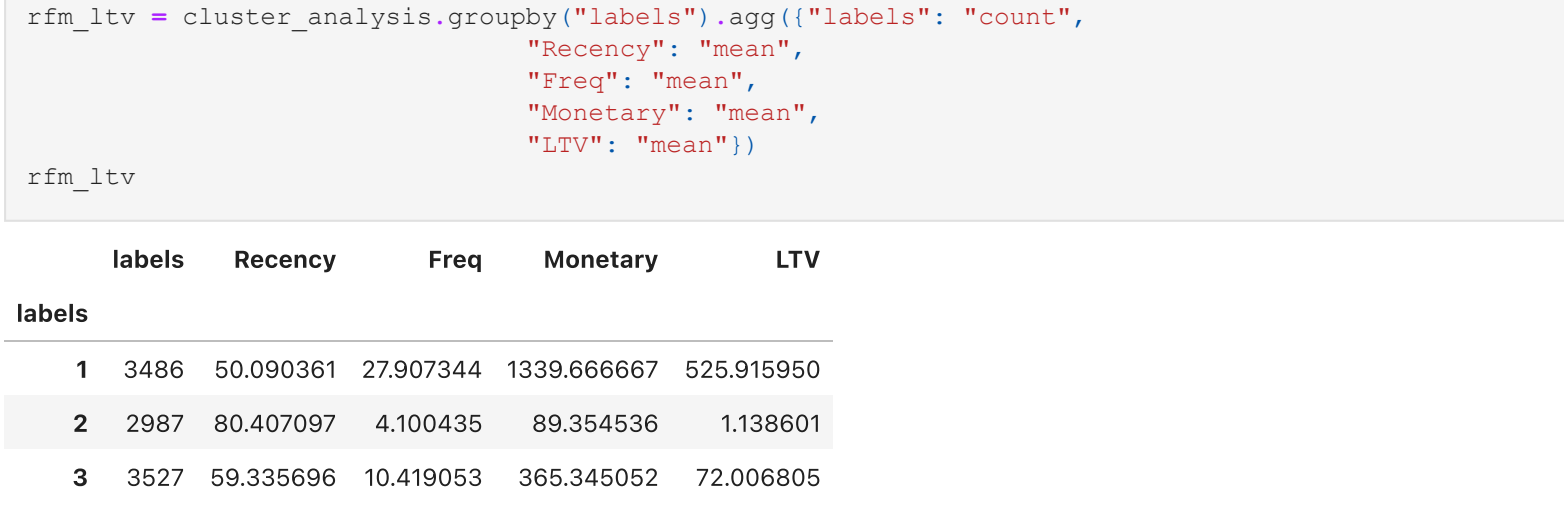
	labels	Age	Income	Edu	WebPurchase	WebVisit	Child
labels							
1	3486	66.690763	98942.249570	16.833333	22.808090	3.110442	0.336776
2	2987	28.707399	40448.831938	15.643120	56.823234	6.685973	0.866756
3	3527	45.659200	66149.769492	17.574142	49.481713	6.053870	0.933088

```
In [48]: # 1: young, lowest lifetime value, purchase web  
# 2: really old, no websearch, almost no webpurchase, og, highest income, CASH COW (LTV)  
# 3: 40-50, 2nd highest income - mittel zwischen den bei  
# 1 and 3 have kids at home
```

```
In [49]: # scaling data for visualization  
scaled_w_clusters = pd.concat((df_scaled, pd.Series(km_labels, name='labels')), axis=1)  
scaled_w_clusters["labels"] = scaled_w_clusters.labels + 1
```

```
In [50]: # scaling data for visualization  
socio_dem = scaled_w_clusters.groupby("labels").agg(  
    "Age": "mean",  
    "Income": "mean",  
    "Edu": "mean")
```

```
In [51]: def line_chart(df):  
    plt.subplots(figsize=(15,8))  
  
    sns.color_palette("Paired")  
    sns.lineplot(data=df, dashes=False)  
  
    plt.xlabel("Features", fontweight='bold')  
    plt.ylabel("Value", fontweight='bold')  
  
    plt.title("lineplot")  
    plt.legend(framealpha=1, frameon=True)  
    plt.show()  
  
def barplot(df):  
    df.T.plot(kind="bar", figsize=(15,8))  
    plt.show()
```



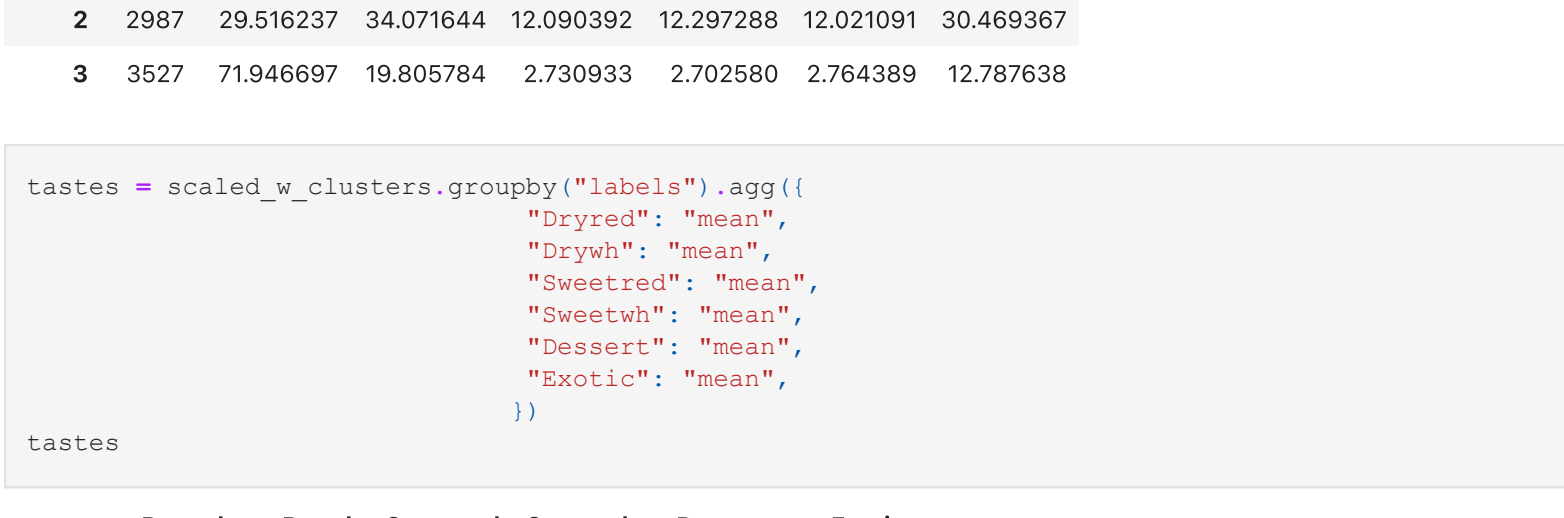
4.4. Profitability RFM Perspective

```
In [54]: rfm_ltv = cluster_analysis.groupby("labels").agg({"labels": "count",  
        "Recency": "mean",  
        "Freq": "mean",  
        "Monetary": "mean",  
        "LTV": "mean"})  
rfm_ltv
```

	labels	Recency	Freq	Monetary	LTV
labels					
1	3486	50.090361	27.907344	1339.666667	525.915950
2	2987	80.407097	4.100435	89.354536	1.138601
3	3527	59.335696	10.419053	365.345052	72.006805

```
In [55]: rfm = scaled_w_clusters.groupby("labels").agg(  
    "Recency": "mean",  
    "Freq": "mean",  
    "Monetary": "mean",  
    "LTV": "mean"  
)  
rfm
```

	Recency	Freq	Monetary	LTV
labels				
1	0.091239	0.489224	0.437842	0.957499
2	0.146461	0.056372	0.027365	0.036799
3	0.108080	0.171256	0.117973	0.126971



```
In [58]: # 1: less recent, not frequent, no money  
# 2: most recent, most money, most frequent, highest LTV ergo CASH COW  
# 3: second recent everywhere
```

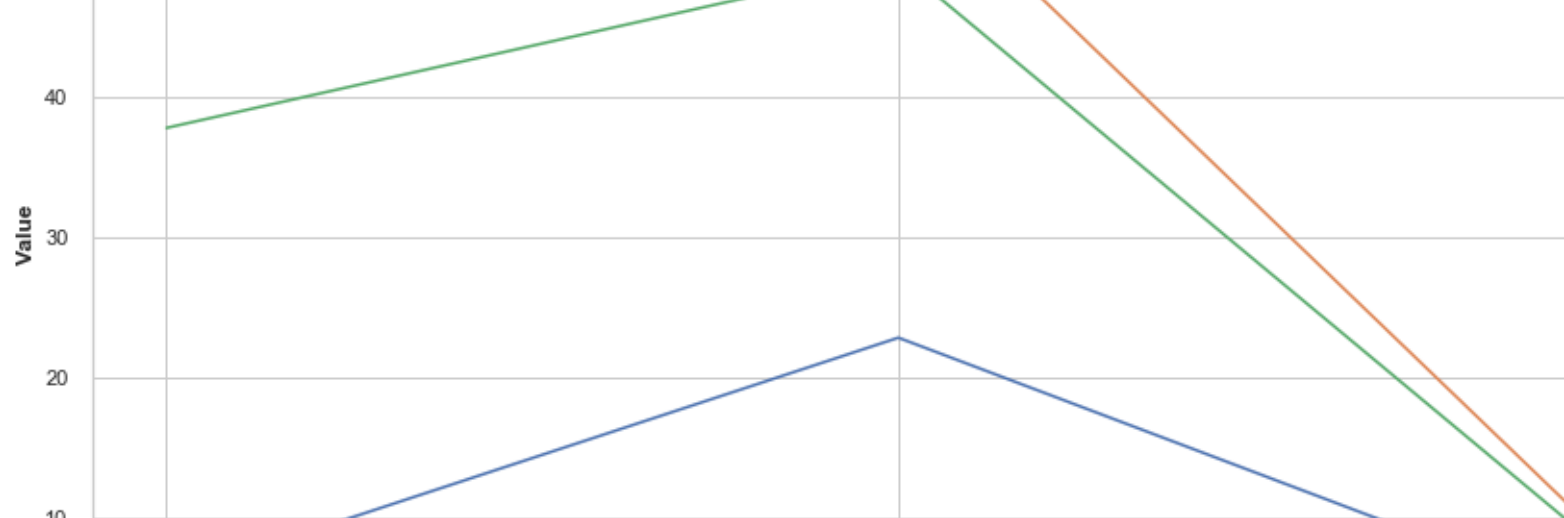
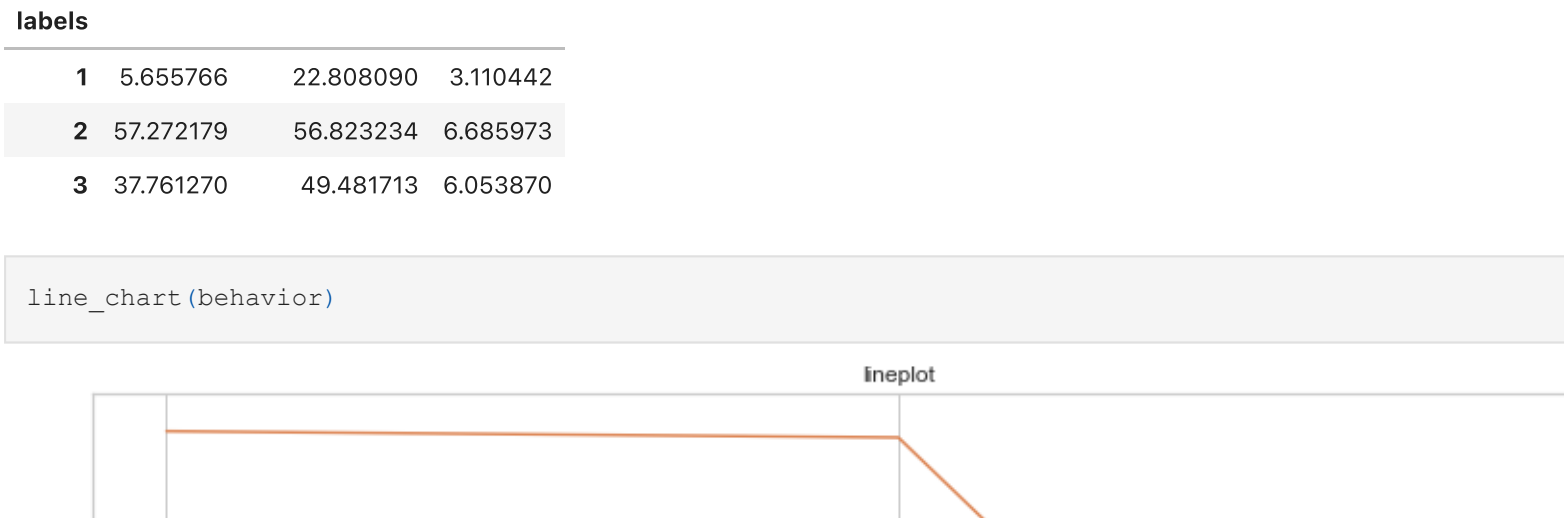
4.5. Taste Perspective

```
In [59]: taste = cluster_analysis.groupby("labels").agg({"labels": "count",  
        "Dryed": "mean",  
        "Drywh": "mean",  
        "Sweetest": "mean",  
        "Sweetwh": "mean",  
        "Dessert": "mean",  
        "Exotic": "mean",  
        "Exotic": "mean",  
        })  
taste
```

	labels	Dryed	Drywh	Sweetest	Sweetwh	Dessert	Exotic
labels							
1	3486	46.444636	32.583477	7.113884	7.009180	6.832186	8.419966
2	2987	29.516237	34.071644	12.090392	12.297288	12.021091	30.469367
3	3527	71.946697	19.805784	2.730933	2.702580	2.764389	12.787638

```
In [60]: tastes = scaled_w_clusters.groupby("labels").agg(  
    "Dryed": "mean",  
    "Drywh": "mean",  
    "Sweetest": "mean",  
    "Sweetwh": "mean",  
    "Dessert": "mean",  
    "Exotic": "mean",  
    "Exotic": "mean",  
    )  
tastes
```

	Dryed	Drywh	Sweetest	Sweetwh	Dessert	Exotic
labels						
1	0.463721	0.432650	0.094852	0.113051	0.088730	0.087708
2	0.290982	0.453036	0.161205	0.198343	0.156118	0.317389
3	0.723946	0.257613	0.036412	0.043590	0.035901	0.133205



```
In [63]: # 1: discount purchase, sweeted -> sweet they like sweet + so krasses zeugs  
# 2: dryed too  
# 3: like dryed a lot, also like to purchase wegg discounted
```

4.6. Buying Behavior Perspective

```
In [64]: behavior = cluster_analysis.groupby("labels").agg(  
    "Perdeal": "mean",  
    "WebPurchase": "mean",  
    "WebVisit": "mean",  
    "WebVisit": "mean",  
    )  
behavior
```

	Perdeal	WebPurchase	WebVisit
labels			
1	5.655766	22.808090	3.110442
2	57.272179	56.823234	6.685973
3	37.761270	49.481713	6.053870



```
In [67]: # 1: young buyers taht purchase online and like to shop on discount  
# 2: dont shop online  
# 3: same as cluster 1
```