

Clustering Crypto

```
In [1]: # Inicial Imports
import pandas as pd
import numpy as np
from os import path
import plotly.express as px
from sklearn.preprocessing import StandardScaler, MinMaxScaler
from sklearn.decomposition import PCA
from sklearn.cluster import KMeans
```

Deliverable 1: Preprocessing the Data for PCA

```
In [7]: # Load the crypto_data.csv dataset.
crypto_df = pd.read_csv("Resources/crypto_data.csv", index_col=0)
```

```
Out[7]:
```

CoinName	Algorithm	IsTrading	ProofType	TotalCoinsMined	TotalCoinSupply
42	42 Coin	Scrypt	PoW/PoS	4.199995e+01	42
365	365Coin	X11	PoW/PoS	NaN	230000000
404	404Coin	Scrypt	PoW/PoS	1.055185e+09	532000000
611	SixEleven	SHA-256	PoW	NaN	611000
808	808	SHA-256	PoW/PoS	0.000000e+00	0
1337	EliteCoin	X13	PoW/PoS	2.927942e+10	314159265359
2015	2015 coin	X11	PoW/PoS	NaN	0
BTC	Bitcoin	SHA-256	PoW	1.792718e+07	21000000
ETH	Ethereum	Ethash	PoW	1.076842e+08	0
LTC	Litecoin	Scrypt	PoW	6.303924e+07	84000000

```
In [8]: crypto_df.shape
```

```
Out[8]: (1252, 6)
```

```
In [9]: # Keep all the cryptocurrencies that are being traded.
crypto_df_traded = crypto_df[crypto_df['IsTrading'] == True]
```

```
Out[9]:
```

CoinName	Algorithm	IsTrading	ProofType	TotalCoinsMined	TotalCoinSupply
42	42 Coin	Scrypt	PoW/PoS	4.199995e+01	42
365	365Coin	X11	PoW/PoS	NaN	230000000
404	404Coin	Scrypt	PoW/PoS	1.055185e+09	532000000
611	SixEleven	SHA-256	PoW	NaN	611000
808	808	SHA-256	PoW/PoS	0.000000e+00	0

```
In [10]: crypto_df_traded.shape
```

```
Out[10]: (1144, 6)
```

```
In [11]: # Keep all the cryptocurrencies that have a working algorithm.
crypto_df_traded_work_algo = crypto_df_traded[crypto_df_traded['Algorithm'].notnull()]
```

```
In [12]: crypto_df_traded_work_algo
```

```
Out[12]:
```

CoinName	Algorithm	IsTrading	ProofType	TotalCoinsMined	TotalCoinSupply
42	42 Coin	Scrypt	PoW/PoS	4.199995e+01	42
365	365Coin	X11	PoW/PoS	NaN	230000000
404	404Coin	Scrypt	PoW/PoS	1.055185e+09	532000000
611	SixEleven	SHA-256	PoW	NaN	611000
808	808	SHA-256	PoW/PoS	0.000000e+00	0
...
SERO	Super Zero	Ethash	PoW	NaN	100000000
UOS	UOS	SHA-256	DPoS	NaN	100000000
BDX	Beldex	CryptoNight	PoW	9.802226e+08	1400222610
ZEN	Horizen	Equihash	PoW	7.296538e+06	21000000
XBC	BitcoinPlus	Scrypt	PoS	1.283270e+05	1000000

```
1144 rows x 6 columns
```

```
In [14]: # Remove the "isTrading" column.
crypto_df_drop_IsTrading = crypto_df_traded_work_algo.drop(['IsTrading'], axis=1)
crypto_df_drop_IsTrading
```

```
Out[14]:
```

CoinName	Algorithm	ProofType	TotalCoinsMined	TotalCoinSupply
42	42 Coin	Scrypt	PoW/PoS	4.199995e+01
365	365Coin	X11	PoW/PoS	NaN
404	404Coin	Scrypt	PoW/PoS	1.055185e+09
611	SixEleven	SHA-256	PoW	NaN
808	808	SHA-256	PoW/PoS	0.000000e+00
...
SERO	Super Zero	Ethash	PoW	NaN
UOS	UOS	SHA-256	DPoS	NaN
BDX	Beldex	CryptoNight	PoW	9.802226e+08
ZEN	Horizen	Equihash	PoW	7.296538e+06
XBC	BitcoinPlus	Scrypt	PoS	1.283270e+05

```
1144 rows x 5 columns
```

```
In [15]: crypto_df_drop_null = crypto_df_drop_IsTrading.dropna()
```

```
Out[15]: crypto_df_drop_null
```

```
Out[15]:
```

CoinName	Algorithm	ProofType	TotalCoinsMined	TotalCoinSupply
42	42 Coin	Scrypt	PoW/PoS	4.199995e+01
404	404Coin	Scrypt	PoW/PoS	1.055185e+09
1337	EliteCoin	X13	PoW/PoS	2.927942e+10
BTC	Bitcoin	SHA-256	PoW	1.792718e+07
ETH	Ethereum	Ethash	PoW	1.076842e+08
...
ZEPH	ZEPHYR	SHA-256	DPoS	2.000000e+09
GAP	Gapcoin	Scrypt	PoW/PoS	1.493105e+07
BDX	Beldex	CryptoNight	PoW	9.802226e+08
ZEN	Horizen	Equihash	PoW	7.296538e+06
XBC	BitcoinPlus	Scrypt	PoS	1.283270e+05

```
685 rows x 5 columns
```

```
In [16]: # Keep the rows where coins are mined.
crypto_df_coins_mined = crypto_df_drop_null[crypto_df_drop_null['TotalCoinsMined'] >= 1]
crypto_df_coins_mined
```

```
Out[16]:
```

CoinName	Algorithm	ProofType	TotalCoinsMined	TotalCoinSupply
42	42 Coin	Scrypt	PoW/PoS	4.199995e+01
404	404Coin	Scrypt	PoW/PoS	1.055185e+09
1337	EliteCoin	X13	PoW/PoS	2.927942e+10
BTC	Bitcoin	SHA-256	PoW	1.792718e+07
ETH	Ethereum	Ethash	PoW	1.076842e+08
...
ZEPH	ZEPHYR	SHA-256	DPoS	2.000000e+09
GAP	Gapcoin	Scrypt	PoW/PoS	1.493105e+07
BDX	Beldex	CryptoNight	PoW	9.802226e+08
ZEN	Horizen	Equihash	PoW	7.296538e+06
XBC	BitcoinPlus	Scrypt	PoS	1.283270e+05

```
532 rows x 5 columns
```

```
In [17]: # Create a new DataFrame that holds only the cryptocurrencies names.
crypto_df_names = pd.DataFrame(crypto_df_coins_mined['CoinName'])
crypto_df_names
```

```
Out[17]:
```

CoinName
42
404
1337
BTC
ETH
ZEPH
GAP
BDX
ZEN
XBC

```
532 rows x 1 columns
```

```
In [18]: # Drop the 'CoinName' column since it's not going to be used on the clustering algorithm.
crypto_df_drop_CoinName = crypto_df_coins_mined.drop(['CoinName'], axis=1)
crypto_df_drop_CoinName
```

```
Out[18]:
```

Algorithm	ProofType	TotalCoinsMined	TotalCoinSupply
42	Scrypt	PoW/PoS	4.199995e+01
404	Scrypt	PoW/PoS	1.055185e+09
1337	X13	PoW/PoS	2.927942e+10
BTC	SHA-256	PoW	1.792718e+07
ETH	Ethash	PoW	1.076842e+08
...
ZEPH	ZEPHYR	SHA-256	DPoS
GAP	Gapcoin	Scrypt	PoW/PoS
BDX	Beldex	CryptoNight	PoW
ZEN	Horizen	Equihash	PoW
XBC	BitcoinPlus	Scrypt	PoS

```
532 rows x 4 columns
```

```
In [19]: # Use get_dummies() to create variables for text features.
X = pd.get_dummies(crypto_df_drop_CoinName, columns=['Algorithm','ProofType'])
```

```
Out[19]:
```

TotalCoinsMined	TotalCoinSupply	Algorithm_1GB	AES_Pattern_Search	Algorithm_536	Algorithm_Argon2d	Algorithm_BLAKE256	Algorithm_BlaKE
4.199995e+01	42	0	0	0	0	0	0
1.055185e+09	532000000	0	0	0	0	0	0
2.927942e+10	314159265359	0	0	0	0	0	0
1.792718e+07	21000000	0	0	0	0	0	0
1.076842e+08	0	0	0	0	0	0	0
2.000000e+09	2000000000	0	0	0	0	0	0
1.493105e+07	250000000	0	0	0	0	0	0
9.802226e+08	1400222610	0	0	0	0	0	0
7.296538e+06	21000000	0	0	0	0	0	0
1.283270e+05	1000000	0	0	0	0	0	0

```
685 rows x 9 columns
```

```
In [20]: # Standardize the data with StandardScaler().
X_scaled = StandardScaler().fit_transform(X)
```

```
Out[20]: array([-0.11710817, -0.152873, -0.0433963, ..., -0.0433963,
```

```
[-0.03936955, -0.145009, -0.0433963, ..., -0.0433963,
```

```
[ 0.52494561, 4.48942416, -0.0433963, ..., -0.0433963,
```

```
..., -0.0433963, -0.0433963, ..., -0.0433963,
```

```
[-0.11694817, -0.15285522, -0.0433963, ..., -0.0433963,
```

```
, -0.0433963, ..., -0.0433963, ..., -0.0433963])
```

```
In [21]: # Create a new DataFrame that has the scaled data with the clustered_df DataFrame index.
```

```
plot_df = pd.DataFrame(X_scaled, columns=["TotalCoinSupply", "TotalCoinsMined"], index=clustered_df.index)
```

```
Out[21]:
```

TotalCoinSupply	TotalCoinsMined	CoinName	Class
4.199995e+01	42	42 Coin	1
1.055185e+09	532000000	404Coin	1
2.927942e+10	314159265359	1337	1
1.792718e+07	21000000	BTC	0
1.076842e+08	0	ETH	0
2.000000e+09	42	ZEPH	1
1.493105e+07	25000000	GAP	1
9.802226e+08	1400222610	BDX	0
7.296538e+06	21000000	ZEN	0
1.283270e+05	1000000	XBC	0

```
532 rows x 4 columns
```

```
In [22]: # Create a new DataFrame that has the scaled data with the clustered_df DataFrame index.
```

```
plot_df["Class"] = clustered_df["Class"]
```

```
plot_df.head(10)
```

```
Out[22]:
```

TotalCoinSupply	TotalCoinsMined	CoinName	Class
4.199995e+01	42	42 Coin	1
1.055185e+09	404Coin	404Coin	1
2.927942e+10	1337	X13	1</