

ITXXX-LAB_ASSIGNMENT1

Akshara

211AI012

Dataset:

WHO life_expectancy data was chosen the dataset consists of 22 columns. I chose 'status' column as label (Developed or developing). This dataset consists of data (like life expectancy, schooling, gdp, per capita income etc) of 193 countries over the period of 2000 to 2015. Using status column as label and the rest 21 as features we can transform this into a binary classification problem where we can classify whether a country is developed or developing for the upcoming years using its features data.

Link: <https://www.kaggle.com/datasets/kumarajarshi/life-expectancy-who>

Google colab notebook:

code can be found here –

<https://colab.research.google.com/drive/1mq9E1rm8wf5eY-kljplvpSaHSy4Du9mT?usp=sharing>

Tasks:

A. Visualize the data and explain why PCA or SVD must be used

In the dataset taken there are 21 features. The training samples size is of only 2350 rows this may lead to the model to overfit (curse of dimensionality). moreover there can be redundant features which depend on each other. This will lead to a poor model.

1.

`<ipython input 10-c186e71d2f1:1> FutureWarning: The default value of numeric_only in DataFrame.corr is deprecated. In a future version, it will default to False. Select only valid columns or specify the value of numeric_only to silence this warning.`

`fig, (ax1, ax2) = plt.subplots(2, 1)`

	Year	Life expectancy	Adult mortality	Infant deaths	Alcohol	percentage expenditure	Hepatitis B	Measles	MM	under-five deaths	Polio	Total expenditure	diphtheria	HIV/AIDS	GDP	population	thinness 16-19 years	thinness 5-9 years	income composition of resources	schooling
Year	1.000000	0.170033	-0.878052	-0.937415	-0.052960	0.031400	0.194333	-0.042403	0.108074	-0.042037	0.004158	0.009740	0.134337	-0.130741	0.101620	0.010060	-0.010776	-0.050020	0.243468	0.209460
Life expectancy	0.170033	1.000000	-0.090359	-0.180557	0.404077	0.301864	0.267062	-0.157580	0.167894	-0.222520	0.460556	0.218066	0.479495	-0.505058	0.461455	-0.021538	-0.477183	-0.471984	0.724778	0.751915
Adult Mortality	-0.878052	-0.090359	1.000000	0.078756	-0.195048	-0.240360	-0.162476	0.031978	-0.307817	0.004148	-0.274823	-0.115201	-0.275131	0.523621	-0.296849	-0.013647	0.303004	0.308457	-0.407628	-0.454617
Infant deaths	-0.937415	-0.180557	0.078756	1.000000	-0.115048	-0.085612	-0.223566	0.001128	-0.227279	0.006629	-0.176689	-0.120616	-0.175171	0.025231	-0.004427	0.556801	0.465711	0.471306	-0.145139	-0.193720
Alcohol	-0.052960	0.404077	-0.195048	-0.115048	1.000000	0.341205	0.007549	-0.051027	0.330408	-0.112379	0.221734	0.206042	0.222020	-0.048845	0.354712	-0.035252	-0.426795	-0.471814	0.450040	0.547378
percentage expenditure	0.031400	0.301864	-0.240360	-0.085612	0.341205	1.000000	0.016274	-0.056096	0.228700	-0.067852	0.147259	0.174420	0.143624	-0.007857	0.000373	-0.025602	-0.251369	-0.252985	0.301992	0.309687
Hepatitis B	0.194333	0.256762	-0.162476	-0.223566	0.007549	0.016274	1.000000	-0.126520	0.150300	-0.233126	0.408171	0.056200	0.011495	-0.112675	0.003303	-0.123321	-0.120429	-0.124980	0.195940	0.231117
Measles	-0.042403	-0.157580	0.031978	0.001128	-0.051027	-0.056096	-0.126520	1.000000	-0.175077	0.007859	-0.138166	-0.106241	-0.141882	0.030099	-0.076466	0.205968	0.224008	0.221072	-0.129568	-0.137225
MM	0.108074	0.004158	-0.307817	-0.227279	0.330408	0.228700	0.150300	-0.175077	1.000000	-0.257080	0.254589	0.242593	0.263147	-0.242517	0.301957	-0.072391	-0.320525	-0.338911	0.500774	0.546861
under-five deaths	-0.042037	-0.222520	0.004148	0.006629	-0.112379	-0.067852	-0.233126	0.007859	-0.257080	1.000000	-0.100729	-0.130165	-0.155568	0.030042	-0.112681	0.544423	0.407780	0.472283	-0.183305	-0.200373
Polio	0.004158	0.460556	-0.274823	-0.176689	0.221734	0.147259	0.408171	-0.138166	0.254589	-0.100729	1.000000	0.137330	0.073593	-0.155660	0.219376	-0.030540	-0.222862	-0.222862	0.301070	0.417860
Total expenditure	0.009740	0.218066	-0.115201	-0.120616	0.206042	0.174420	0.056200	-0.102401	0.242503	-0.138165	0.137330	1.000000	0.152754	-0.001380	0.138364	-0.079602	-0.277181	-0.283774	0.166862	0.246384
Diphtheria	0.134337	0.479495	-0.275131	0.175171	0.222020	0.143624	0.011495	-0.141882	0.201147	-0.155568	0.073593	0.152754	1.000000	-0.164880	0.209666	-0.030444	-0.229518	-0.222743	0.401456	0.425352
HIV/AIDS	-0.130741	-0.505058	0.523621	0.025231	-0.048845	-0.007857	-0.112675	0.030099	-0.243717	0.030002	-0.155660	-0.001380	-0.164880	1.000000	-0.136491	-0.027854	0.204064	0.207283	-0.249518	-0.220409
GDP	0.101620	0.461455	-0.296849	-0.104427	0.354712	0.000373	0.003303	-0.076466	0.201507	-0.112681	0.219376	0.138364	0.209666	-0.136491	1.000000	-0.030778	-0.205687	-0.206039	0.400341	0.448273
Population	0.010060	-0.021538	-0.303004	0.556801	-0.035252	-0.025602	-0.123321	0.205968	-0.072391	-0.544423	-0.030540	-0.079602	-0.028444	-0.027854	-0.030778	1.000000	0.253044	0.251483	-0.006735	-0.031668
thinness 16-19 years	-0.047628	-0.477183	0.302904	0.465711	-0.426795	-0.251369	-0.126429	0.224008	-0.532025	0.407780	-0.221823	-0.277181	-0.229518	0.204064	-0.205687	0.253044	1.000000	0.939182	-0.425429	-0.471652
thinness 5-9 years	-0.050020	-0.471984	0.308457	0.471306	-0.471414	-0.252985	-0.124980	0.221072	-0.538911	0.472283	-0.222862	-0.283774	-0.222743	0.207283	-0.206039	0.251483	0.939182	1.000000	-0.411653	-0.460632
income composition of resources	0.243468	0.724778	-0.407628	-0.145139	0.450040	0.301992	0.199549	-0.129568	0.000774	-0.183305	0.301070	0.166862	0.401456	-0.249519	0.400341	-0.008735	-0.422429	-0.411653	1.000000	0.000000
Schooling	0.209460	0.751915	-0.454617	-0.193720	0.547378	0.309687	0.231117	-0.137225	0.546861	-0.200373	0.417860	0.246384	0.425352	-0.220429	0.448273	-0.031668	-0.471652	-0.460632	0.000000	1.000000

This is the correlation matrix of the data we can see that some features are highly correlated with each other. Examples : schooling and life expectancy(0.75), infant deaths and underfive deaths(0.99), polio and diphtheria(0.67), income composition and schooling(0.8) etc. there are a lot of features that have correlation >0.4 or <-0.4 between them. This explains that we have remove correlation in features in our data. So we use SVD or PCA. SVD or PCA can transform the original features into a new set of uncorrelated variables, known as principal components. This can be beneficial because

it **reduces the impact of multicollinearity**, making the model more robust and interpretable.

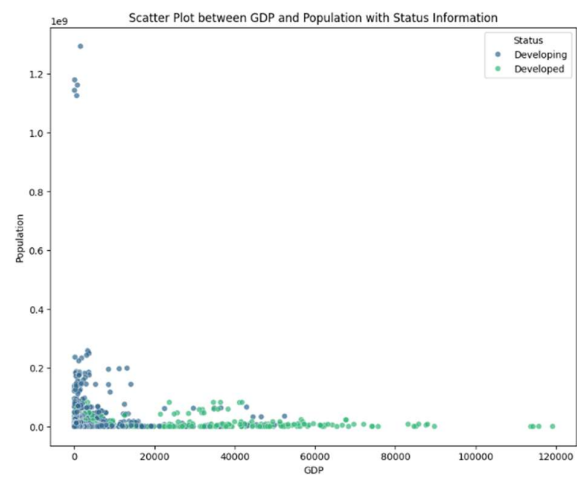
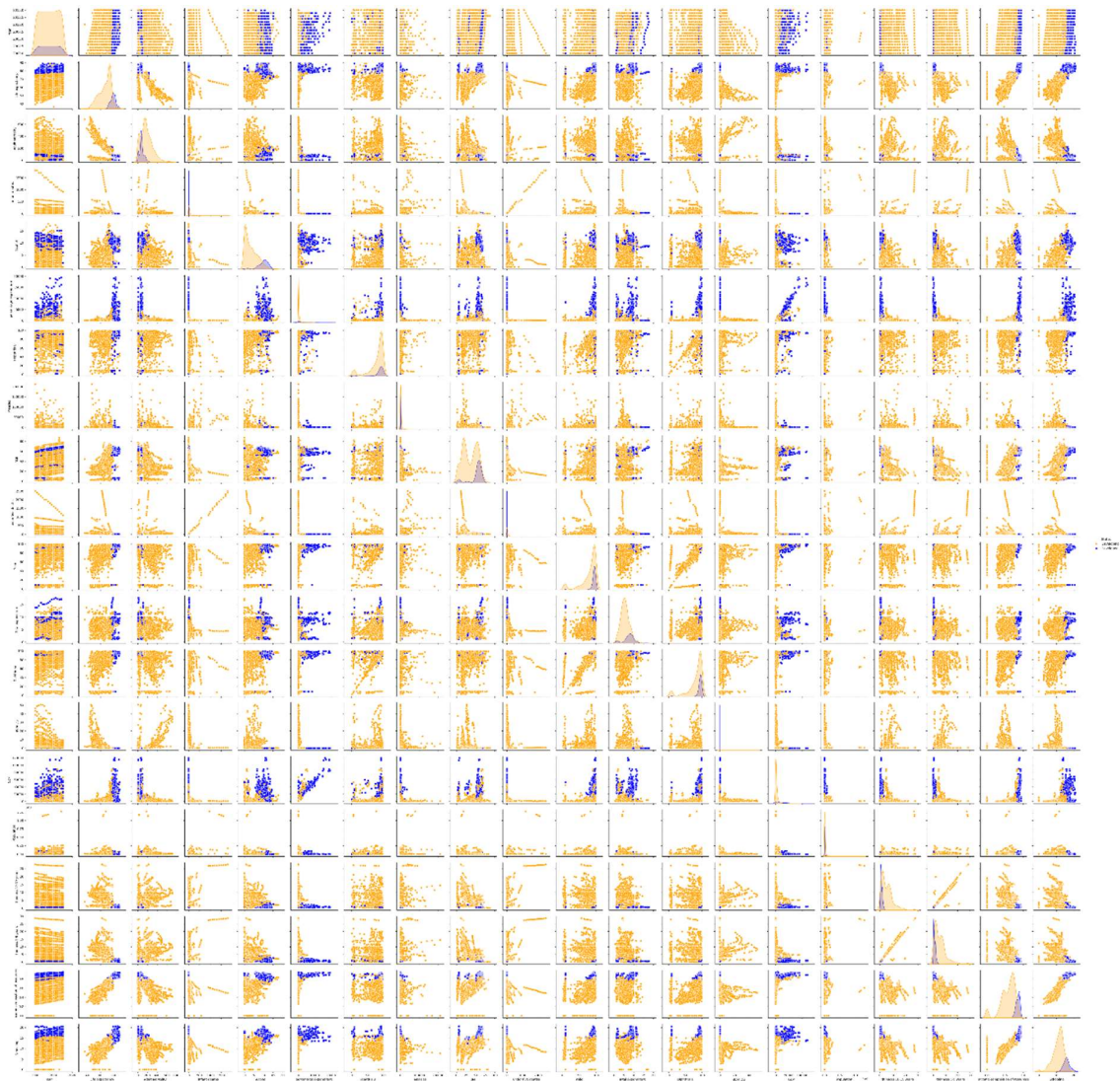
2.

Covariance matrix: (python Input:11-ff0f37937af5): FutureWarning: The default value of numeric_only in DataFrame.cov is deprecated. In a future version, it will default to False. Select only valid columns or specify the value of numeric_only to silence this warning. covariance_matrix = df[['col1', 'col2']]																				
	Year	Life expectancy	Adult mortality	Infant deaths	Alcohol	percentage expenditure	Hepatitis B	Malaria	meas	under-five deaths	Polio	total expenditure	diphtheria	measles	gdp	Population	thinness 10-19 years	thinness 5-9 years	income composition of resources	Schooling
Year	2.128753e+01	7.486219e+00	-4.530899e+01	-2.835746e+01	-8.319021e-01	2.879802e+02	1.147744e+01	-4.364552e+03	1.007855e+01	-3.179481e+01	1.017423e+01	9.821331e+01	1.489486e+01	-3.272873e+08	6.876510e+03	4.771411e+06	-8.762649e-01	-1.699350e+00	6.236081	3.244610e+00
Life expectancy	7.486219e+00	8.679485e+01	-8.243083e+02	-2.711084e+02	1.566435e+01	7.248051e+03	5.412226e+01	-1.723845e+04	1.975874e+02	5.405781e+02	1.035024e+02	5.161025e+09	1.078214e+02	-2.895623e+01	6.357257e+04	-1.260191e+07	-2.061529e+01	-2.817786e+01	1.431375	2.355088e+01
Adult mortality	-4.530899e+01	-8.243083e+02	1.548823e+04	1.156183e+03	9.839165e+01	-4.989727e+04	-4.884433e+02	4.450702e+04	-9.568895e+02	1.888458e+02	-7.980218e+02	-3.931276e+01	-8.088108e+02	3.310213e+02	-6.530286e+05	-1.074518e+05	1.689953e+02	1.722356e+02	-11.871124	-1.870545e+02
Infant deaths	-2.835746e+01	-2.711084e+02	1.156183e+03	1.368886e+04	5.027989e+01	-2.988884e+04	-3.868975e+02	6.778735e+05	-5.402218e+02	1.885706e+04	-4.739808e+02	-3.881213e+01	-4.114747e+02	1.538885e+01	-1.959717e+05	4.474351e+09	2.441691e+02	2.545259e+02	-3.678165	-7.910588e+01
Alcohol	1.566435e+01	1.566435e+01	8.839165e+01	1.164285e+01	1.620591e+02	2.811414e+03	1.888889e+01	-2.054987e+03	2.678386e+01	-7.425908e+01	2.988374e+01	3.088816e+00	2.138245e+01	-1.838387e+00	2.985325e+04	-8.988531e+06	-7.731891e+09	7.677487e+00	8.379128	7.189164e+00
percentage expenditure	2.879802e+02	7.248051e+03	-4.989727e+04	-2.988884e+04	2.831143e+03	3.951895e+06	4.485416e+02	-1.298152e+06	9.159427e+03	-3.822058e+04	6.877933e+03	8.964324e+02	6.750915e+03	-8.877866e+02	2.738311e+07	-3.380358e+09	-2.226007e+03	-2.734889e+03	164.211694	2.666523e+03
Hepatitis B	1.147744e+01	5.412226e+01	-4.884433e+02	-5.068751e+02	8.838323e+00	6.480514e+02	6.285957e+02	-2.881926e+04	7.520127e+01	-4.223259e+02	2.654561e+02	3.638453e+00	3.288785e+02	-1.446184e+01	2.881837e+04	-2.143913e+03	-1.292949e+01	-1.358915e+01	6.969997	1.643195e+01
Malaria	-4.364552e+03	-1.723845e+04	4.450702e+04	6.778735e+05	-2.454987e+03	-1.298153e+06	-2.881926e+04	1.314803e+06	-4.967268e+04	8.348334e+05	-3.889537e+04	-3.118151e+03	-3.870728e+04	1.789181e+03	-1.208813e+07	1.673944e+11	1.146315e+04	1.148208e+04	-268.290238	-4.859832e+03
meas	1.007855e+01	1.007855e+01	8.888889e+02	-5.422226e+02	2.879486e+01	8.192421e+03	7.828101e+01	-4.687289e+04	4.617823e+02	7.688139e+02	1.332898e+02	1.910224e+01	1.343213e+02	-2.483234e+01	1.632626e+04	-8.838424e+07	-4.176661e+01	-8.784889e+01	2.160844	3.593358e+01
under-five deaths	3.174871e+01	3.405781e+02	1.888458e+02	1.885706e+04	7.428389e+01	-2.988884e+04	2.232538e+02	8.434304e+05	-7.683178e+02	2.937477e+04	-7.185389e+02	4.361628e+01	-7.488237e+02	3.188823e+01	-3.781758e+05	6.941725e+09	3.183774e+02	3.432356e+02	6.628414	1.145498e+02
Polio	1.017423e+01	1.035024e+02	-7.980218e+02	-4.739808e+02	2.988374e+01	4.877863e+03	2.654561e+02	-3.668317e+04	1.332898e+02	7.185389e+02	5.488733e+02	7.838817e+00	3.742538e+02	-1.893481e+01	7.883131e+04	-5.688020e+07	-2.889378e+01	-2.340278e+01	1.814838	3.134777e+01
Total expenditure	9.821331e+01	5.161025e+09	-3.931276e+01	-3.881121e+01	3.088816e+00	8.964324e+02	3.638453e+00	-3.118151e+03	1.219221e+01	5.351835e+01	7.888171e+00	6.241681e+00	8.977258e+00	-1.820789e+02	4.827722e+03	-1.188041e+07	-3.093672e+00	-3.232271e+00	8.883395	1.948959e+00
diphtheria	1.489486e+01	1.078214e+02	-2.895623e+01	-4.914747e+02	2.138245e+01	6.750915e+03	3.288785e+02	-3.878728e+02	1.343278e+02	-7.488237e+02	3.742538e+02	8.977258e+00	5.624891e+02	-1.989886e+01	6.728158e+04	-4.188462e+07	-2.401881e+01	-2.378178e+01	1.820803	3.205273e+01
measles	-3.272873e+08	-2.895623e+01	3.318823e+02	1.510885e+01	-1.058751e+00	-8.877866e+02	-1.446184e+01	1.789181e+03	-2.493284e+01	3.188893e+01	-1.963485e+01	-1.825766e+02	-1.999886e+01	-1.058789e+04	-8.589218e+05	4.603716e+09	4.778178e+00	4.287385	-3.836737e+00	
GDP	6.876510e+03	6.357257e+04	-5.338286e+05	-1.957971e+05	2.065328e+04	2.738831e+07	2.681837e+04	-1.208813e+07	8.826256e+04	-2.751755e+05	7.883131e+04	4.827722e+03	6.727815e+04	-1.858796e+04	2.838277e+08	-2.348734e+10	-1.841628e+04	-1.911338e+04	1481.319128	2.178588e+04
Population	4.771411e+06	-1.260191e+07	-1.074518e+08	4.474351e+09	-8.988531e+06	-3.388385e+09	-2.143913e+08	1.873944e+11	-8.838424e+07	5.947125e+09	-5.688020e+07	-1.188043e+07	-4.188462e+07	-8.585218e+06	-3.457346e+10	3.722476e+15	7.782216e+07	7.356578e+07	-12882.747229	4.653202e+00
thinness 10-19 years	8.762649e-01	-2.061529e+01	1.688883e+02	2.441691e+02	-7.731891e+00	-2.226007e+03	-1.292949e+01	1.148308e+04	-4.735861e+01	3.332744e+02	-2.888783e+01	-3.893672e+00	-2.481081e+00	4.683718e+00	-1.841628e+04	7.782216e+07	1.953812e+01	1.671643e+01	-4.382794	-8.868058e+00
thinness 5-9 years	-1.699350e+00	-2.817786e+01	1.722925e+02	2.538228e+02	-7.077487e+00	-2.278488e+03	-1.358915e+01	1.148308e+04	-4.870408e+01	3.432526e+02	-2.348259e+01	-3.232271e+00	-2.377879e+01	4.778178e+00	-1.811386e+04	7.758615e+07	1.871643e+01	2.838021e+01	-4.388888	4.838818e+00
income composition of resources	2.388886e-01	1.431375e+00	-1.187112e+01	-3.679185e+00	3.791275e-01	1.642117e+02	8.888872e-01	-2.882802e+02	2.100044e+00	5.628414e+00	1.814838e+00	8.304484e+02	1.820063e+00	-2.728731e+01	1.481138e+03	-1.126527e+05	-3.827935e+01	-3.888838e-01	8.844488	5.644489e-01
Schooling	3.244610e+00	2.355088e+01	-1.879458e+01	-7.815889e+01	7.216184e+00	2.688323e+03	1.643195e+01	-4.883832e+03	3.330385e+01	1.148308e+02	3.134777e+01	1.848828e+00	3.268273e+01	-3.887231e+08	2.178588e+04	-4.632831e+08	-8.888888e+00	-4.838818e+00	8.844447	1.128234e+01

This is a covariance matrix. This also demonstrates that some features like life expectancy and income composition, infant deaths and under five deaths are covariant. So we need to reduce redundancy and dimensionality of the data for which PCA or SVD can be used. SVD and PCA can help us identify and retain the most important features, **reducing the dimensionality** while preserving most of the variance in the data. This can help remove noise and focus on the essential information.

3.

As we can see below we have scatter pair plot of every column with every other column (zoom in for better view).the blue dots represent developed and orange samples represent developing. We can see some features have linear(or nearly linear) plots like infant deaths and under five deaths, gdp and percentage expenditure, Thinness(5-9) and thinness(10-19). This demonstrates the collinearity in our data. PCA can help address collinearity issues by transforming them into uncorrelated principal components. Moreover Reducing the dimensionality of data a few principal components can facilitate better visualization, allowing us to explore and understand the structure of the data. This can be particularly helpful when dealing with high-dimensional datasets.

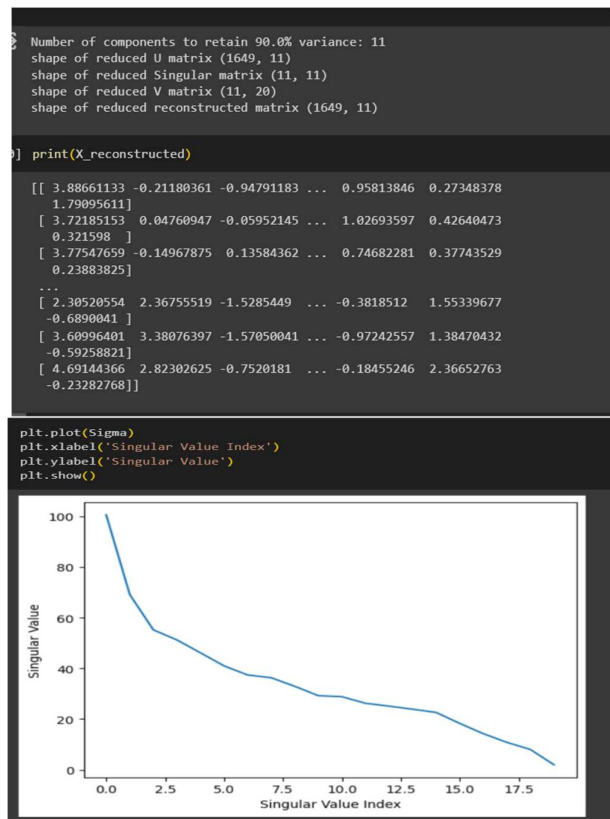


B. Implement SVD and PCA logic on your own and find the appropriate k-dimensions for the data

1. SVD

- The `svd_decomposition` function takes the matrix `X_std` and performs Singular Value Decomposition (SVD). The result includes left singular vectors (`U`), singular values (`Sigma`), and right singular vectors (`Vt`). This function is implemented using block power iteration method. Block Power Iteration iteratively refines an initial approximation to the left singular vectors of the matrix through the computation of matrix products and QR decompositions. It is particularly useful when dealing with large matrices where computing the full SVD might be computationally expensive.
- Then we decide `k` based on desired variance to retain. We get `k=11` for 90% variance. we calculate cumulative explained variance from the singular values and determines the number of components (`k`) needed to retain a specified amount of variance (`desired_variance`).
- We retain only the top `k` components of the SVD results: left singular vectors (`U_k`), a diagonal matrix of top `k` singular values (`Sigma_k`), and the top `k` right singular vectors (`Vt_k`).
- It reconstructs the data using the reduced SVD components: `U_k`, `Sigma_k`, and `Vt_k`. we can do this either by `U_k.dot(Sigma_k)` or `(X_std).dot(Vt_k.T)`

Output:



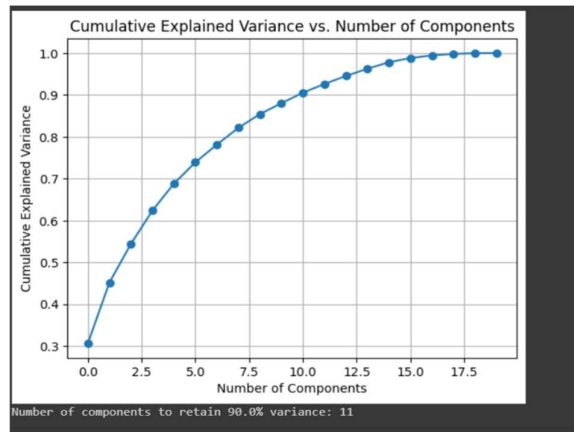
When we plot singular values we see that only about 10 or so singular value indexes have significantly high singular values the rest have very low. This demonstrates that we

can summarise most of the variance in our data by choosing an appropriate k (here we chose 11)

2. PCA

- a. First we standardise the data and compute the covariance matrix from the standardized data. The covariance matrix represents the relationships between different features in the dataset
- b. Then we compute the eigenvalues and corresponding eigenvectors of the covariance matrix. Eigenvectors represent the directions of maximum variance, and eigenvalues indicate the magnitude of the variance in those directions.
- c. We sort the eigenvalues in descending order and rearranges the corresponding eigenvectors accordingly. This step is crucial as it allows identifying the principal components in order of significance.
- d. Next we the cumulative explained variance, which represents the proportion of total variance explained by each principal component and the sum of all previous components.
- e. We plot the cumulative explained variance against the number of components. This visualization helps in determining how many components are needed to retain a certain percentage of the total variance.
- f. We Set a threshold (desired_variance=90%) for the cumulative explained variance and determines the minimum number of components (k) needed to achieve at least that much variance.
- g. We get k=11
- h. Then we chose top 11 eigenvectors based on their eigen values and project it onto initial X_Std data to get the final transformed data

Output:




```

Dimensions of standardized_data: (1649, 20)
Dimensions of top_eigenvectors: (20, 11)
Dimensions of transformed_data: (1649, 11)

transformed_data
array([[ -3.88661133,  0.21180361, -0.9478715 , ...,  0.95814995,
        -0.30265613,  1.7862947  ],
       [ -3.72185153, -0.04760947, -0.05959068, ...,  1.02694563,
        -0.4315925 ,  0.31454981 ],
       [ -3.77547659,  0.14967875,  0.13575916, ...,  0.74683278,
        -0.38127993,  0.23260972 ],
       ...,
       [ -2.30520554, -2.36755519, -1.52908612, ..., -0.38185435,
        -1.54195562, -0.71400818 ],
       [ -3.60996401, -3.38076397, -1.5714142 , ..., -0.9724252 ,
        -1.37485687, -0.6149685  ],
       [ -4.69144366, -2.82302625, -0.75307605, ..., -0.18453999,
        -2.36241681, -0.2712815  ]])

```

We can see that we got the same transformed data using either SVD or PCA.

C. Visualize the data (t-sne plot) after applying SVD and PCA

A t-SNE model is created with two components, indicating a two-dimensional visualization.

Fit and Transform Data Using t-SNE:

The t-SNE model is fitted and used to transform the data into a lower-dimensional space.

Map 'status' Labels to Numeric Values:

The original categorical labels ('Developed', 'Developing') are mapped to numeric values (0, 1) for color-coding.

Create Scatter Plot:

A scatter plot is generated with points in the t-SNE-transformed space. Points are colored based on the numeric labels and have reduced transparency.

Plot Configuration:

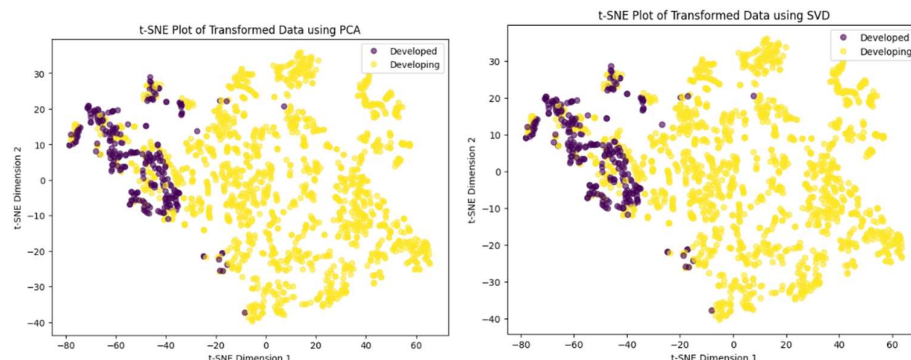
The plot is configured with a title, and axis labels for the two dimensions.

Add Legend:

A legend is added to the plot, indicating the correspondence between numeric labels and the original categorical labels.

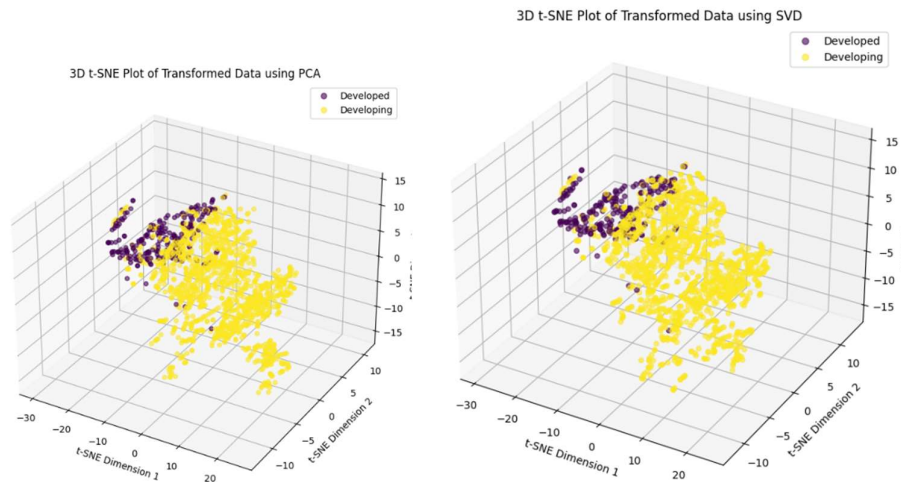
Show Plot:

The final t-SNE plot is displayed for visualization.



Similarly a 3D tsne scatter plot is plotted.

Here we use top 3 principal componenets to plot the graph using developed and developing sample data points.



D. State your conclusions as to how SVD and PCA have helped here

- SVD and PCA are effective techniques for reducing the dimensionality of the dataset while preserving essential information. This is particularly useful in our dataset with a large number of features in the dataset.
- The cumulative explained variance plot obtained through PCA assists in determining the number of principal components needed to retain a certain percentage of the total variance. This aids in balancing dimensionality reduction with the retention of information.
- The transformed data obtained after applying SVD and PCA can be used for further analysis and interpretation. We can build better performing model using the transformed data
- Initially when we plot pair plot between features we could see that feature vs feature plot was very difficult to fit a binary classification model. But the TSNE plot using top two principal components has a clean boundary and its easy to build a classifier model on this transformed data using principal components.
- These principal components can be used as features to build better and efficient models.

using original data					using transformed data				
Accuracy: 0.7998					Accuracy: 0.9273				
Confusion Matrix:					Confusion Matrix:				
[[288 31]					[[73 25]				
[296 1098]]					[23 539]]				
Classification Report:					Classification Report:				
	precision	recall	f1-score	support		precision	recall	f1-score	support
Developed	0.41	0.87	0.56	239	Developed	0.76	0.74	0.75	98
Developing	0.97	0.79	0.87	1394	Developing	0.96	0.96	0.96	562
accuracy			0.80	1633	accuracy			0.93	660
macro avg	0.69	0.83	0.72	1633	macro avg	0.86	0.85	0.85	660
weighted avg	0.89	0.88	0.82	1633	weighted avg	0.93	0.93	0.93	660

- I have built an SVM classifier to test the accuracy and we see that our model performs better when trained with transformed data with reduced dimensions.
- In summary, the application of SVD and PCA to the WHO dataset has provided a means of dimensionality reduction, visualization, and insights into the relationships between features and life expectancy categories.