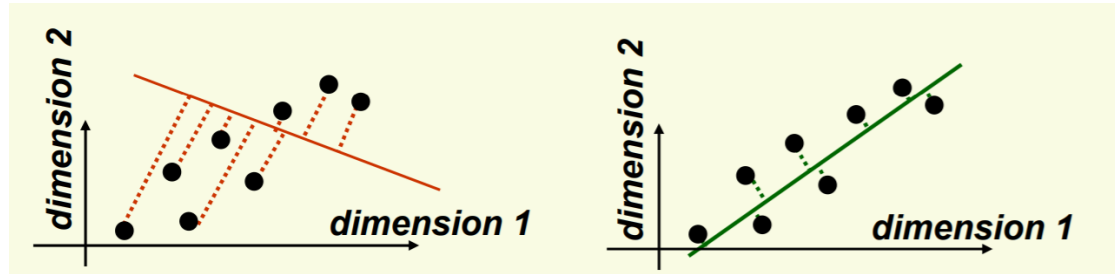


Principal Component Analysis Concepts

Principal Component Analysis

1. Main idea: seek most accurate data representation in a lower dimensional space
1. Example in 2-D, project data to 1-D subspace (a line) with minimal projection error

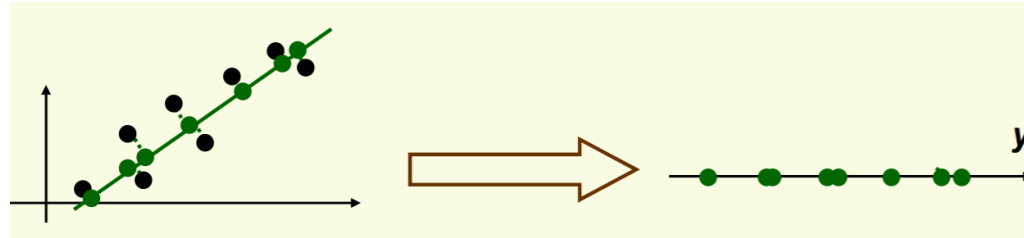


1. In both the pictures above, the data points (black dots) are projected to one line but the second line is closer to the actual points (less projection errors) than first one
1. Notice that the good line to use for projection lies in the direction of largest variance

Ref: http://www.cs.haifa.ac.il/~rita/uml_course/add_mat/PCA.pdf

Principal Component Analysis

5. After the data is projected on the best line, need to transform the coordinate system to get 1D representation for vector y

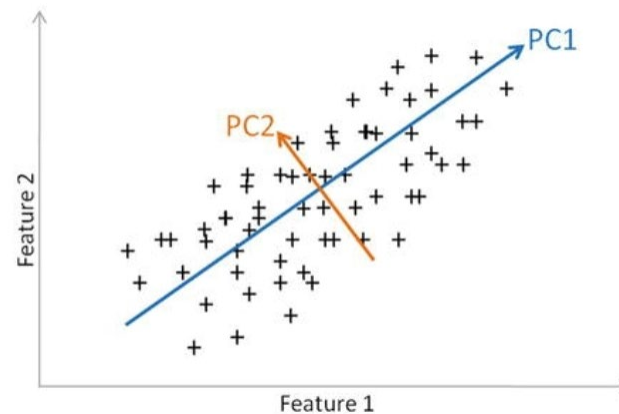


5. Note that new data y has the same variance as old data x in the direction of the green line
5. PCA preserves largest variances in the data

Ref: http://www.cs.haifa.ac.il/~rita/uml_course/add_mat/PCA.pdf

Principal Component Analysis

8. In general PCA on n dimensions will result in another set of new n dimensions. The one which captures maximum variance in the underlying data is the principal component 1, principal component 2 is orthogonal to it
8. Example in 2-D, project data to 1-D subspace (a line) with minimal projection error



Ref: http://www.cs.haifa.ac.il/~rita/uml_course/add_mat/PCA.pdf

Mechanics of Principal Component Analysis

<http://setosa.io/ev/principal-component-analysis/>

Principal Component Analysis steps

1. Begins by standardizing the data. Data on all the dimensions are subtracted from their means to shift the data points to the origin. i.e. the data is centered on the origins
1. Generate the covariance matrix / correlation matrix for all the dimensions
1. Perform eigen decomposition, that is, compute eigen vectors which are the principal components and the corresponding eigen values which are the magnitudes of variance captured
1. Sort the eigen pairs in descending order of eigen values and select the one with the largest value. This is the first principal component that covers the maximum information from the original data

Ref: http://www.cs.haifa.ac.il/~rita/uml_course/add_mat/PCA.pdf

Principal Component Analysis (Performance issues)

1. PCA effectiveness depends upon the scales of the attributes. If attributes have different scales, PCA will pick variable with highest variance rather than picking up attributes based on correlation
1. Changing scales of the variables can change the PCA
1. Interpreting PCA can become challenging due to presence of discrete data
1. Presence of skew in data with long thick tail can impact the effectiveness of the PCA (related to point 1)
1. PCA assumes linear relationship between attributes. It is ineffective when relationships are non linear

Principal Component Analysis steps

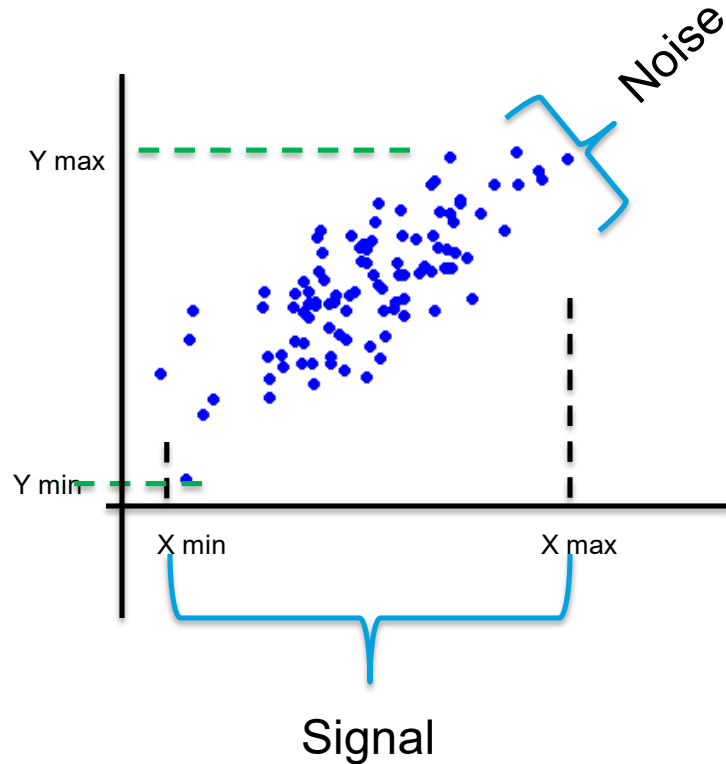
Lab-3 Principal Component Analysis on iris data set

Description — Explore the iris data set and perform PCA

The data set is winequality-red.csv

Sol: PCA-iris.ipynb

Principal Component Analysis (Signal to noise ratio)



```
X_std_df = pd.DataFrame(X_std)
axes = pd.plotting.scatter_matrix(X_std_df)
plt.tight_layout()
```

Signal – all valid values for a variable (shown between max and min values for x axis and y axis). Represents a valid data

Noise – The spread of data points across the best fit line. For a given value of x, there are multiple values of y (some on line and some around the line). This spread is due to random factors

Signal to Noise Ratio – Variance of signal / variance in noise. $\frac{\sigma_{signal}^2}{\sigma_{noise}^2}$

Greater the SNR the better the model will be

Principal Component Covariance Matrix

1. Variance is measured within the dimensions and co-variance is among the dimensions

$$\text{var}(X) = \frac{\sum_{i=1}^n (X_i - \bar{X})(X_i - \bar{X})}{(n-1)}$$

1. Express total variance (variance and cross variance between dimensions as a matrix (variance matrix)

$$\text{cov}(X, Y) = \frac{\sum_{i=1}^n (X_i - \bar{X})(Y_i - \bar{Y})}{(n-1)}$$

1. Covariance matrix is a mathematical representation of the total variance of individual dimension and across dimensions .

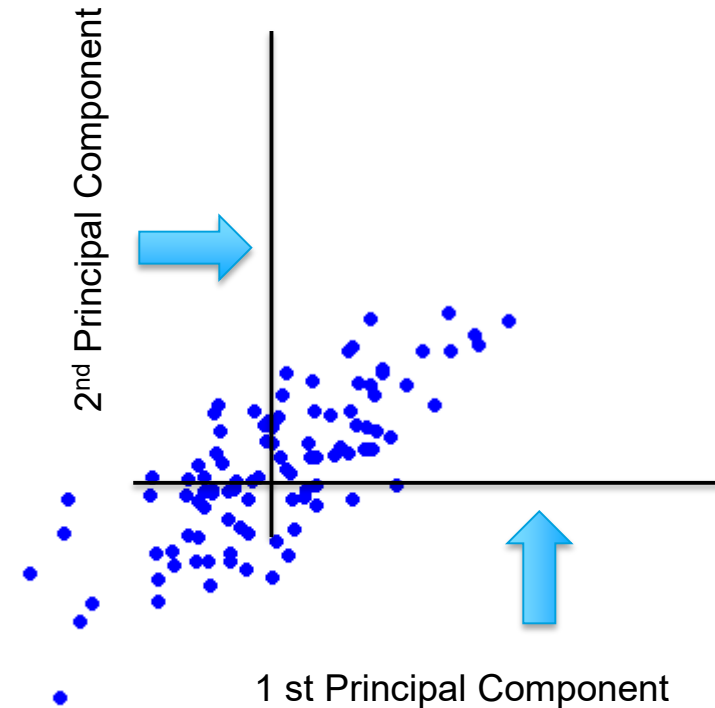
$$C = \begin{bmatrix} \text{cov}(X, X) & \text{cov}(X, Y) & \text{cov}(X, Z) \\ \text{cov}(Y, X) & \text{cov}(Y, Y) & \text{cov}(Y, Z) \\ \text{cov}(Z, X) & \text{cov}(Z, Y) & \text{cov}(Z, Z) \end{bmatrix}$$

Covariance matrix for three dimensions x,y and z

`eig_vals, eig_vecs = np.linalg.eig(cov_matrix)`

Improving SNR through PCA (Scaling the dimensions)

1. The mean is subtracted from all the points on both dimensions i.e. $(x_i - \bar{x})$ and $(y_i - \bar{y})$
1. The dimensions are transformed using algebra into new set of dimensions
1. The transformation is a rotation of axes in mathematical space



```
X_std = StandardScaler().fit_transform(X)
eig_vals, eig_vecs = np.linalg.eig(cov_matrix)
```

PCA (Calculating total variance (covariance and variance)

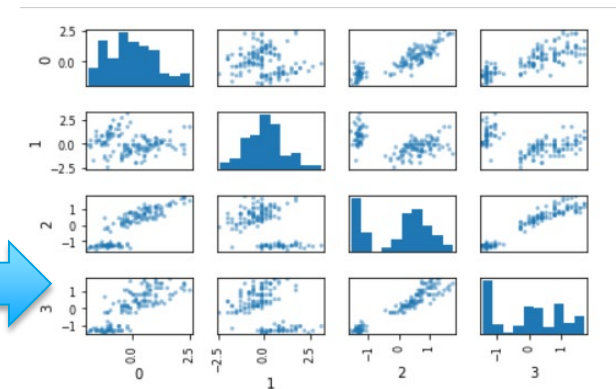
- Multiplying the two matrices produces a matrix of total variance also called covariance matrix (a square and symmetric matrix).

$$\begin{matrix} A \\ \begin{bmatrix} a_{11} & a_{12} & \dots & a_{1n} \\ a_{21} & a_{22} & \dots & a_{2n} \\ \dots & \dots & \dots & \dots \\ a_{m1} & a_{m2} & \dots & a_{mn} \end{bmatrix} \end{matrix} \times \begin{matrix} A^T \\ \begin{bmatrix} a_{11} & a_{21} & \dots & a_{m1} \\ a_{12} & a_{22} & \dots & a_{m2} \\ \dots & \dots & \dots & \dots \\ a_{1n} & a_{2n} & \dots & a_{mn} \end{bmatrix} \end{matrix} = C = \begin{bmatrix} \text{cov}(X,X) & \text{cov}(X,Y) & \text{cov}(X,Z) \\ \text{cov}(Y,X) & \text{cov}(Y,Y) & \text{cov}(Y,Z) \\ \text{cov}(Z,X) & \text{cov}(Z,Y) & \text{cov}(Z,Z) \end{bmatrix}$$



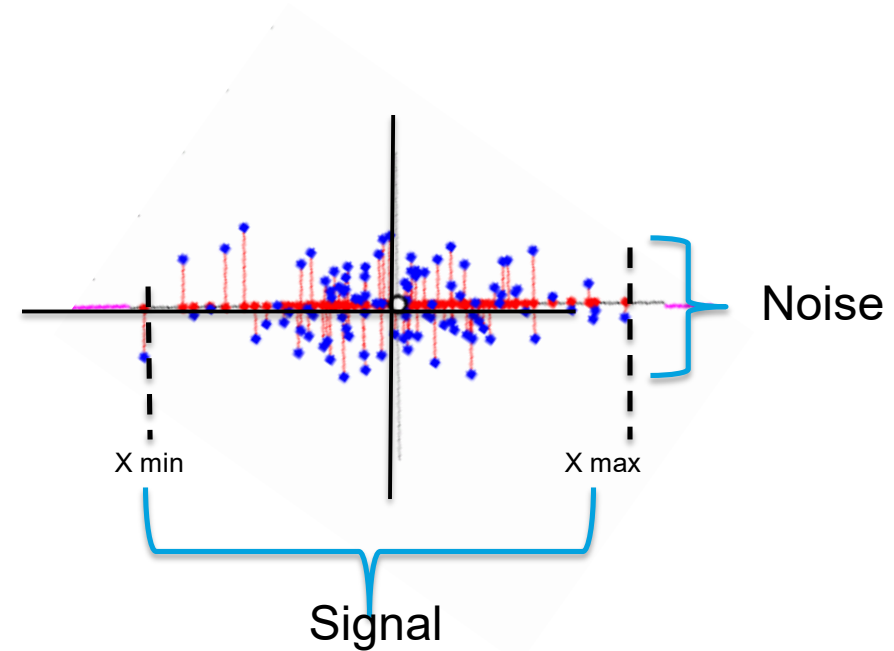
Covariance Matrix

```
%s [[ 1.00671141 -0.11010327  0.87760486  0.82344326]
[-0.11010327  1.00671141 -0.42333835 -0.358937 ]
[ 0.87760486 -0.42333835  1.00671141  0.96921855]
[ 0.82344326 -0.358937    0.96921855  1.00671141]]
```



Improving SNR through PCA (Principal components)

5. The original data points are now represented by the red dots on new dimensions
5. It also introduces error of representation (vertical red lines from the blue dots to corresponding red dots on the new dimension)
5. The axis rotation is done such that the new dimension captures max variance in the data points and also reduces total error of representation



```
print('Eigen Vectors \n%s', eig_vecs)
print('\n Eigen Values \n%s', eig_vals)
```

Properties of principal components and their covariance matrix

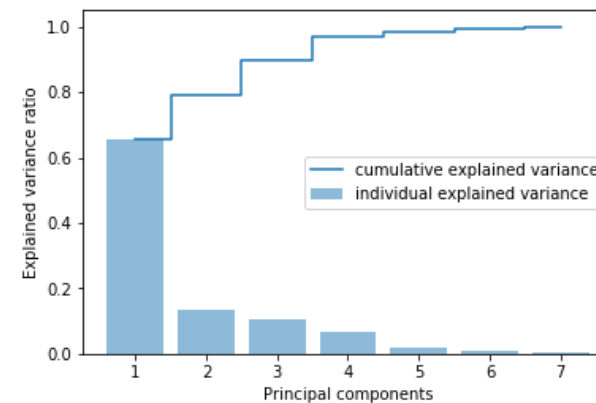
8. Thus to find principal components we need to get the diagonal matrix $\mathbf{B}\mathbf{B}^T$ from the original covariance matrix $\mathbf{A}\mathbf{A}^T$

$$\begin{array}{c}
 \begin{bmatrix} \text{cov}(X,X) & \text{cov}(X,Y) & \text{cov}(X,Z) \\ \text{cov}(Y,X) & \text{cov}(Y,Y) & \text{cov}(Y,Z) \\ \text{cov}(Z,X) & \text{cov}(Z,Y) & \text{cov}(Z,Z) \end{bmatrix} \\
 \mathbf{A}\mathbf{A}^T
 \end{array}
 \quad \longrightarrow \quad
 \begin{array}{c}
 \begin{pmatrix} \sigma_1^2 & 0 & \dots & 0 \\ 0 & \sigma_2^2 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & \sigma_p^2 \end{pmatrix} \\
 \mathbf{B}\mathbf{B}^T
 \end{array}$$

8. For this we have to transform the matrix \mathbf{A} to a new matrix \mathbf{B} such that the covariance matrix of \mathbf{B} ($\mathbf{B}\mathbf{B}^T$), is a **diagonal matrix** (Ref to part 2, bullet 5)

PCA for dimensionality reduction

1. PCA can also be used to reduce dimensions
1. Arrange all eigen vectors along with corresponding eigen values in descending order of eigen values
1. Plot a cumulative eigen_value graph as shown below
1. Eigen vectors with insignificant contribution to total eigen values can be removed from analysis (for e.g. eigen vector 6 and 7 below)



Thanks