

Objective: To develop a mathematical understanding of PCA

Test data

```
In [12]: import numpy as np
from scipy import linalg as LA
X = np.array([
    [0.387, 4878, 5.42],
    [0.723, 12104, 5.25],
    [1, 12756, 5.52],
    [1.524, 6787, 3.94],
])
X = X - np.mean(X, axis=0)
```

Non-linear Iterative Partial Least-Squares (NIPALS) algorithm

Steps to compute PCA using NIPALS algorithm

- Step 1: Initialize an arbitrary column vector \mathbf{t}_a either randomly or by just copying any column of X .
- Step 2: Take every column of X , X_k and regress it onto the \mathbf{t}_a vector and store the regression coefficients as p_{ka} . (Note: This simply means performing an ordinary least squares regression ($y = mx$) with $x = t_a$ and $y = X_k$ with $m = (\mathbf{x}^T \mathbf{x})^{-1} \mathbf{x}^T \mathbf{y}$). In the current notation we get

$$p_{ka} = \frac{\mathbf{t}_a^T X_k}{\mathbf{t}_a^T \mathbf{t}_a}$$

Repeat it for each of the columns of X to get the entire vector \mathbf{p}_k . This is shown in the illustration above where each column from X is regressed, one at a time, on \mathbf{t}_a , to calculate the loading entry, p_{ka}

In practice we don't do this one column at a time; we can regress all columns in X in go:

$$\mathbf{p}_a^T = \frac{1}{\mathbf{t}_a^T \mathbf{t}_a} \cdot \mathbf{t}_a^T X_a$$

where \mathbf{t}_a is an $N \times 1$ column vector, and X_a is an $N \times K$ matrix.

- The loading vector \mathbf{p}_a^T won't have unit length (magnitude) yet. So we simply rescale it to have magnitude of 1.0:

$$\mathbf{p}_a^T = \frac{\mathbf{p}_a^T}{\sqrt{\mathbf{p}_a^T \mathbf{p}_a}}$$

- Step 4: Regress every row in X onto this normalized loadings vector. As illustrated below, in our linear regression the rows in X are our y-variable each time, while the loadings vector is our x-variable. The regression coefficient becomes the score value for that i^{th} row:

$$p_{i,a} = \frac{\mathbf{x}_i^T \mathbf{p}_a}{\mathbf{p}_a^T \mathbf{p}_a}$$

where \mathbf{x}_i^T is a $K \times 1$ column vector. We can combine these N separate least-squares models and calculate them in one go to get the entire vector,

$$\mathbf{t}_a^T = \frac{1}{\mathbf{p}_a^T \mathbf{p}_a} \mathbf{X} \mathbf{p}_a^T$$

where \mathbf{p}_a is a $K \times 1$ column vector.

- Step 5: Continue looping over steps 2,3,4 until the change in vector t_a is below a chosen tolerance
- Step 6: On convergence, the score vector and the loading vectors, \mathbf{t}_a and \mathbf{p}_a are stored as the a^{th} column in matrix \mathbf{T} and \mathbf{P} . We then deflate the \mathbf{X} matrix. This crucial step removes the variability captured in this component (t_a and p_a) from \mathbf{X} :

$$E_a = X_a - t_a p_a^T$$

$$X_{a+1} = E_a$$

For the first component, X_a is just the preprocessed raw data. So we can see that the second component is actually calculated on the residuals E_1 , obtained after extracting the first component. This is called deflation, and nicely shows why each component is orthogonal to the others. Each subsequent component is only seeing variation remaining after removing all the others; there is no possibility that two components can explain the same type of variability. After deflation we go back to step 1 and repeat the entire process for the next component.

IMPLEMENTATION IN PYTHON

```
In [28]: def PCA(X,no_components):
    tol = 0.0000001
    it=1000
    obsCount,varCount = X.shape
    Xa = X - np.mean(X, axis = 0)
    #Xh = X-np.tile(np.mean(X,axis=0).reshape(-1,1).T, obsCount).reshape(4,3)
    T = np.zeros((obsCount,no_components))
    P = np.zeros((varCount,no_components))
    pcvar = np.zeros((varCount,1))
    varTotal = np.sum(np.var(Xa,axis=0,ddof=1))
    currVar = varTotal
    nr=0
    for h in range(no_components):
        th = Xa[:,0].reshape(-1,1)
        ende = False
        while ende != True:
            nr = nr + 1
            ph = np.dot(Xa.T,th)/np.dot(th.T,th)
            ph = ph /np.linalg.norm(ph)
            thnew = np.dot(Xa,ph)/np.dot(ph.T,ph)
            prec = np.dot((thnew-th).T,(thnew-th))
            th = thnew
            if prec <= (tol*tol):
                ende = True
            elif it <=nr:
                ende = True
                print("Iteration stops without convergence")
        Ea = Xa - np.dot(th,ph.T)
```

```

Xa = Ea
T[:,h] = th.flatten()
P[:,h] = ph.flatten()
oldVar = currVar
currVar = np.sum(np.var(Xa,axis=0,ddof=1))
pcvar[h] = (oldVar - currVar) / varTotal
nr = 0
return T,P,pcvar

```

Advantages of the NIPALS algorithm

- The NIPALS algorithm computes one component at a time. The first component computed is equivalent to the t1 and p1 vectors that would have been found from an eigenvalue or singular value decomposition.
- The algorithm can handle missing data in X.
- The algorithm always converges, but the convergence can sometimes be slow.
- It is also known as the Power algorithm to calculate eigenvectors and eigenvalues.
- It works well for very large data sets.
- It is used by most software packages, especially those that handle missing data.
- Of interest: it is well known that Google used this algorithm for the early versions of their search engine, called PageRank148.

```

In [14]: no_components=3
         T,P,pcvar = PCA(X,no_components)
         print("T (Scores)")
         print(T)
         print(" ")
         print("P (Loadings)")
         print(P)
         print(np.sqrt(pcvar)/np.sum(np.sqrt(pcvar)))

T (Scores)
[[-4.25324997e+03 -8.41288672e-01  8.37859036e-03]
 [ 2.97275001e+03 -1.25977272e-01 -1.82476780e-01]
 [ 3.62475003e+03 -1.56843494e-01  1.65224286e-01]
 [-2.34425007e+03  1.12410944e+00  8.87390330e-03]]

P (Loadings)
[[ 1.21901390e-05  5.66460728e-01  8.24088735e-01]
 [ 9.99999997e-01  5.32639787e-05 -5.14047689e-05]
 [ 7.30130279e-05 -8.24088733e-01  5.66460726e-01]]
[[9.99753412e-01]
 [2.10083377e-04]
 [3.65048880e-05]]

```

SVD

```

In [15]: from numpy.linalg import svd
         U, S, PTrans = svd(X, full_matrices=False)
         Sigma = np.diag(S)
         T=np.dot(U,Sigma)
         P=PTrans.T

         print("T (Scores)")
         print(T)
         print(" ")
         print("P (Loadings)")
         print(P)

```

```
print("Sigma (Variance)")
print(S)
```

```
T (Scores)
[[-4.25324997e+03  8.41288672e-01 -8.37858943e-03]
 [ 2.97275001e+03  1.25977271e-01  1.82476780e-01]
 [ 3.62475003e+03  1.56843494e-01 -1.65224286e-01]
 [-2.34425007e+03 -1.12410944e+00 -8.87390454e-03]]

P (Loadings)
[[ 1.21901390e-05 -5.66460727e-01 -8.24088736e-01]
 [ 9.99999997e-01 -5.32639789e-05  5.14047691e-05]
 [ 7.30130279e-05  8.24088734e-01 -5.66460725e-01]]
Sigma (Variance)
[6.74994067e+03  1.41840009e+00  2.46466604e-01]
```

SKLEARN PCA

```
In [16]: from sklearn.decomposition import PCA
pca = PCA()
T=pca.fit_transform(X)
Prans=pca.components_ #eigen vectors.T
latent = pca.explained_variance_
explained = pca.explained_variance_ratio_
P=PTrans.T
S=pca.singular_values_
Sigma=np.diag(S)
print("T (Scores)")
print(T)
print(" ")
print("P (Loadings)")
print(P)
print("Sigma (Variance)")
print(S)
#print(pca.singular_values_/np.sqrt(3))
```

```
T (Scores)
[[ 4.25324997e+03 -8.41288672e-01 -8.37858943e-03]
 [-2.97275001e+03 -1.25977271e-01  1.82476780e-01]
 [-3.62475003e+03 -1.56843494e-01 -1.65224286e-01]
 [ 2.34425007e+03  1.12410944e+00 -8.87390454e-03]]

P (Loadings)
[[ 1.21901390e-05 -5.66460727e-01 -8.24088736e-01]
 [ 9.99999997e-01 -5.32639789e-05  5.14047691e-05]
 [ 7.30130279e-05  8.24088734e-01 -5.66460725e-01]]
Sigma (Variance)
[6.74994067e+03  1.41840009e+00  2.46466604e-01]
```

```
In [17]: pca.explained_variance_ratio_
```

```
Out[17]: array([9.99999955e-01, 4.41567976e-08, 1.33326424e-09])
```

```
In [18]: explained_variance_2 = (S ** 2) / 4
explained_variance_ratio_2 = (explained_variance_2 / explained_variance_2.sum())
print(explained_variance_ratio_2)

[9.99999955e-01  4.41567976e-08  1.33326424e-09]
```

Eigenvalue decomposition approach

Recall that the latent variable directions (the loading vectors) were oriented so that the variance of the scores in that direction were maximal. We can cast this as an optimization problem. For

the first component:

$$\max(\phi) = \mathbf{t}_1^T \mathbf{t}_1 = \mathbf{p}_1^T \mathbf{X}^T \mathbf{X} \mathbf{p}_1$$

such that

$$\mathbf{p}_1^T \mathbf{p}_1 = 1$$

This is equivalent to

$$\max(\phi) = \mathbf{p}_1^T \mathbf{X}^T \mathbf{X} \mathbf{p}_1 - \lambda(\mathbf{p}_1^T \mathbf{p}_1 - 1)$$

because we can move the constraint into the objective function with a Lagrange multiplier, λ . The maximum value must occur when the partial derivatives with respect to \mathbf{p}_1 ,

our search variable, are zero:

$$\frac{\partial \phi}{\partial \mathbf{p}_1} = \frac{\partial(\mathbf{p}_1^T \mathbf{X}^T \mathbf{X} \mathbf{p}_1 - \lambda(\mathbf{p}_1^T \mathbf{p}_1 - 1))}{\partial \mathbf{p}_1} = 0$$

$$2\mathbf{X}^T \mathbf{X} \mathbf{p}_1 - 2\lambda_1 \mathbf{p}_1 = 0$$

$$(\mathbf{X}^T \mathbf{X} - \lambda_1 \mathbf{I}) \mathbf{p}_1 = 0$$

$$\mathbf{X}^T \mathbf{X} \mathbf{p}_1 = \lambda_1 \mathbf{p}_1$$

which is just the eigenvalue equation, indicating that \mathbf{p}_1 is the eigenvector of $\mathbf{X}^T \mathbf{X}$ and λ_1 is the eigenvalue. One can show that $\lambda_1 = \mathbf{t}_1^T \mathbf{t}_1$, which is proportional to the variance of the first component. In a similar manner we can calculate the second eigenvalue, but this time we add the additional constraint that $\mathbf{p}_1 \perp \mathbf{p}_2$. Writing out this objective function and taking partial derivatives leads to showing that

$$\mathbf{X}^T \mathbf{X} \mathbf{p}_2 = \lambda_2 \mathbf{p}_2$$

From this we learn that:

- The loadings are the eigenvectors of $\mathbf{X}^T \mathbf{X}$.
- Sorting the eigenvalues in order from largest to smallest gives the order of the corresponding eigenvectors, the loadings.
- We know from the theory of eigenvalues that if there are distinct eigenvalues, then their eigenvectors are linearly independent (orthogonal).
- We also know the eigenvalues of $\mathbf{X}^T \mathbf{X}$ must be real values and positive; this matches with the interpretation that the eigenvalues are proportional to the variance of each score vector.
- Also, the sum of the eigenvalues must add up to sum of the diagonal entries of $\mathbf{X}^T \mathbf{X}$, which represents of the total variance of the \mathbf{X} matrix, if all eigenvectors are extracted. So plotting the eigenvalues is equivalent to showing the proportion of variance explained in \mathbf{X} by each component. This is not necessarily a good way to judge the number of components to use, but it is a rough guide: use a Pareto plot of the eigenvalues (though in the context of eigenvalue problems, this plot is called a scree plot).

```
In [19]: cov = np.cov(X, rowvar = False)
evals , P = LA.eigh(cov)
idx = np.argsort(evals)[::-1]
P = P[:,idx]
evals = evals[idx]
T = np.dot(X, P)
Sigma=LA.norm(T,axis=0)
print("T (Scores)")
print(T)
print("P (Loadings)")
print(P)
print("Sigma (Variance)")
print(Sigma)

T (Scores)
[[ 4.25324997e+03  8.41288672e-01  8.37858943e-03]
 [-2.97275001e+03  1.25977271e-01 -1.82476780e-01]
 [-3.62475003e+03  1.56843494e-01  1.65224286e-01]
 [ 2.34425007e+03 -1.12410944e+00  8.87390454e-03]]
P (Loadings)
[[-1.21901390e-05 -5.66460727e-01  8.24088736e-01]
 [-9.99999997e-01 -5.32639789e-05 -5.14047691e-05]
 [-7.30130279e-05  8.24088734e-01  5.66460725e-01]]
Sigma (Variance)
[6.74994067e+03  1.41840009e+00  2.46466604e-01]
```

Task 1: Test if the loading vectors are orthogonal and orthonormal or not

```
In [20]: dot1 = np.dot(P[:,0], P[:,1])
dot2 = np.dot(P[:,1], P[:,2])
dot3 = np.dot(P[:,0], P[:,1])
print(dot1)

dot4 = P.T.dot(P)
print(dot4)

-4.472333961502706e-19
[[ 1.00000000e+00 -4.48550678e-19 -3.82279445e-19]
 [-4.48550678e-19  1.00000000e+00  2.10601467e-16]
 [-3.82279445e-19  2.10601467e-16  1.00000000e+00]]
```

Orthonormal

Task 2: Test if the scores vectors are orthogonal and orthonormal or not

```
In [30]: dotT = T.T.dot(T)
print(dotT)

[[ 4.55616990e+07 -1.95286340e-11 -1.55750968e-11]
 [-1.95286340e-11  2.01185881e+00  4.06619183e-15]
 [-1.55750968e-11  4.06619183e-15  6.07457869e-02]]
```

Orthogonal, not orthonormal.

Task 3: Add more columns to the original data matrix by:

- Make some of the columns to be the linear combination of others
- Duplicate some columns
- Add noise as some columns

Then apply PCA to the dataset and report your findings here

```
In [29]: noise = np.random.normal(scale = 1, size = (X.shape[0],1))

X_alt = np.hstack([X, X[:,2][:,None], 2*X[:,0][:,None] - X[:,1][:,None], noise])
X_alt = X_alt - np.mean(X_alt, axis=0)
print(np.mean(X_alt, axis=0))

T_alt, P_alt, var_alt = PCA(X_alt, 6)
print(np.diag(T_alt.T.dot(T_alt)))

print("Length of P-s", np.diag(P_alt.T.dot(P_alt)))

print("P", P_alt.T.dot(P_alt))

[2.77555756e-17 0.00000000e+00 0.00000000e+00 0.00000000e+00
 1.13686838e-13 5.55111512e-17]
[9.11211784e+07 4.77474771e+00 9.47991363e-01 1.33303002e-25
 1.20403253e-27 1.46941337e-31]
Length of P-s [1. 1. 1. 1. 1. 1.]
P [[ 1.00000000e+00 -1.07542143e-13 -3.60930266e-14 -8.55433489e-01
    5.17912270e-01 -8.69564978e-05]
 [-1.07542143e-13  1.00000000e+00  6.6633950e-19  5.91024167e-05
    7.82195217e-05  9.52692294e-01]
 [-3.60930266e-14  6.6633950e-19  1.00000000e+00 -1.18996428e-05
    -9.28560522e-04  7.73849921e-02]
 [-8.55433489e-01  5.91024167e-05 -1.18996428e-05  1.00000000e+00
    1.50858562e-16 -5.22843266e-14]
 [ 5.17912270e-01  7.82195217e-05 -9.28560522e-04  1.50858562e-16
    1.00000000e+00  2.76166024e-15]
 [-8.69564978e-05  9.52692294e-01  7.73849921e-02 -5.22843266e-14
    2.76166024e-15  1.00000000e+00]]
```

No longer orthogonal. Assume that when adding dependent columns, x-basis is no longer orthogonal, which means that the P-space is also no longer orthogonal.