

Calcul Numeric

Curs 13

2020

Anca Ignat

Principal Component Analysis (PCA)

“a way of identifying patterns in data, and expressing the data in such a way to highlight their similarities and differences”

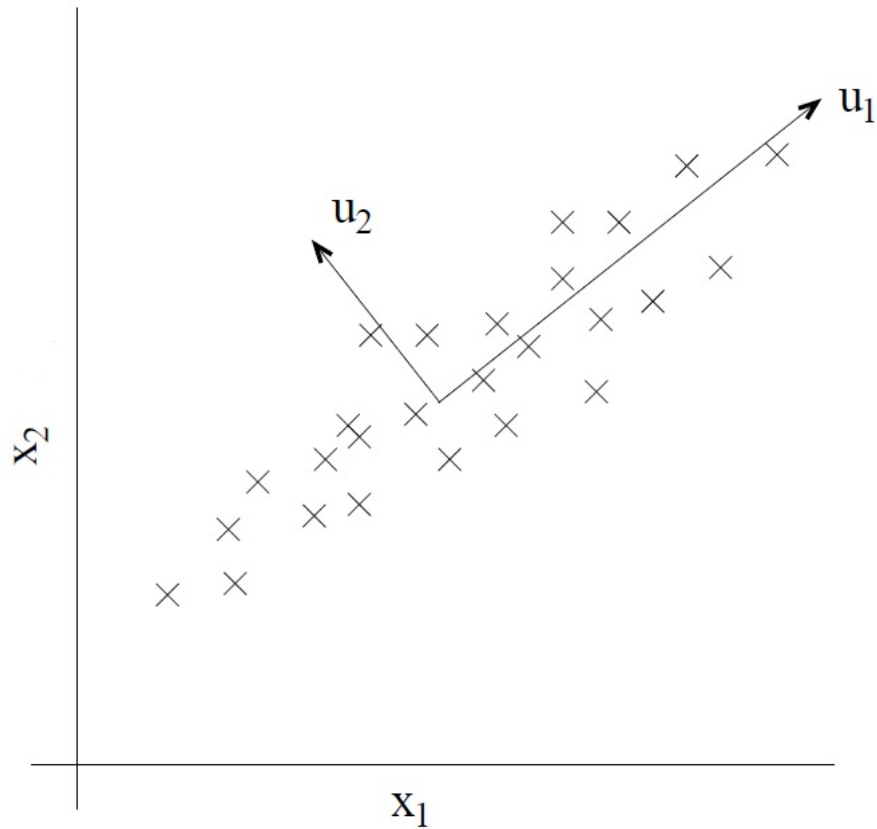
$X = \{\mathbf{x}^1, \mathbf{x}^2, \dots, \mathbf{x}^p\}, \mathbf{x}^k \in \mathbb{R}^n$ - a set of observations

$X \subseteq Y \subset \mathbb{R}^n, \dim Y = k \ll n$ - dimension reduction

(removing redundancies in data)

- some components of the vectors \mathbf{x}^k may be linearly dependant one with the other

(\mathbf{x}_i =grades of a student, \mathbf{x}_j =credits obtained by the student)



u_1 – the principal direction of the data
How to automatically compute u_1 ?

Pre-processing of data

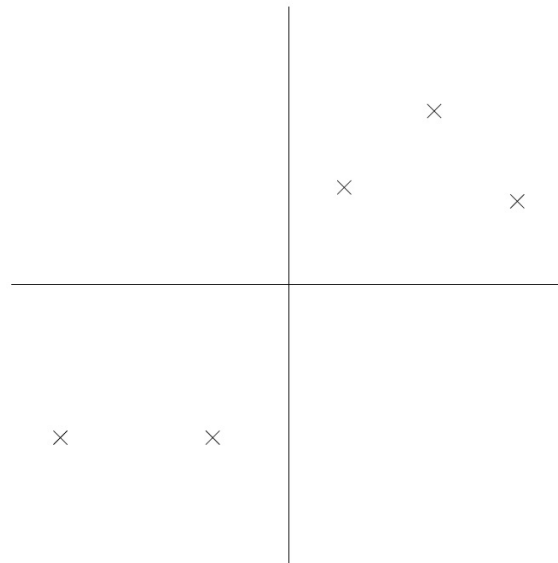
1. Compute the mean value $\mu = \frac{1}{p} \sum_{k=1}^p x^k$
2. Subtract mean from data: replace x^k with $x^k - \mu$
3. Compute $\sigma_i^2 = \frac{1}{p} \sum_{k=1}^p (x_i^k)^2$, $i = 1, \dots, n$
4. Replace each x_i^k with $\frac{1}{\sigma_i} x_i^k$

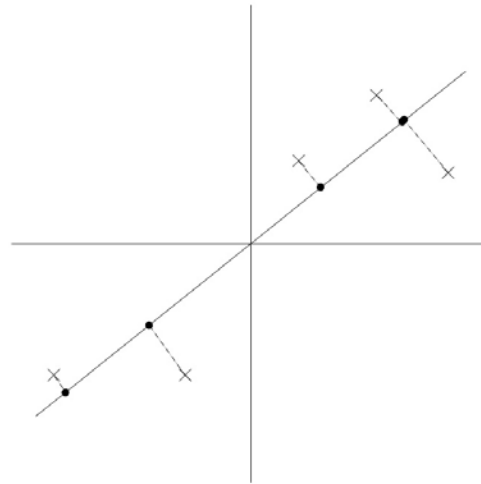
(1)+(2) – creates data with zero mean

(3)+(4) – scale normalization (when the components were computed with different measure units, marks \leftrightarrow credits)

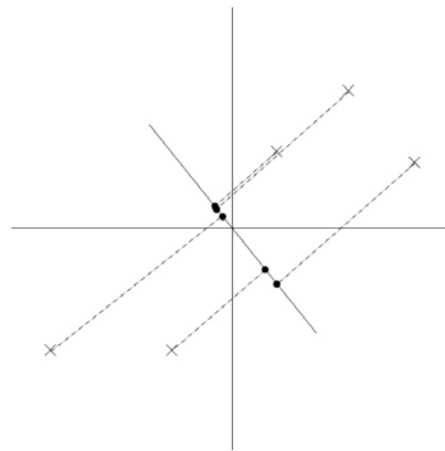
Compute the “major axis of variation” – \mathbf{u} – the direction that describes the data positioning/ alignment.

Find the unit vector \mathbf{u} such that projecting the data onto the direction given by \mathbf{u} one obtains new data with maximized variance.





“good” direction



“bad” direction

The length of the projection of vector \mathbf{x} onto direction \mathbf{u} is given by the inner product $\mathbf{x}^T \mathbf{u} = (\mathbf{x}, \mathbf{u})$.

To maximize the variance of projections one needs to solve the optimization problem:

$$\begin{aligned} \max \{ & \frac{1}{p} \sum_{k=1}^p \left((\mathbf{x}^k)^T \mathbf{u} \right)^2 = \frac{1}{p} \sum_{k=1}^p \mathbf{u}^T \mathbf{x}^k (\mathbf{x}^k)^T \mathbf{u} = \\ & \mathbf{u}^T \left(\frac{1}{p} \sum_{k=1}^p \mathbf{x}^k (\mathbf{x}^k)^T \right) \mathbf{u} \ ; \ \mathbf{u} \in \mathbb{R}^n, \|\mathbf{u}\|_2 = 1 \} \end{aligned}$$

The solution to this problem is the **unit eigenvector** corresponding to the **biggest eigenvalue** of the **covariance matrix \mathbf{C}** of the dataset (with zero mean!) \mathbf{X} .

$$C = \frac{1}{p} \sum_{k=1}^p \mathbf{x}^k (\mathbf{x}^k)^T$$

Consider the eigenvalues of the covariance matrix in decreasing order:

$$\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_n$$

The covariance matrix being symmetric, one can compute an orthonormal basis of eigenvectors $\mathbf{u}^1, \mathbf{u}^2, \dots, \mathbf{u}^n$ - this are the principal components for dataset \mathbf{X} . One can consider only the first most important r principal components, and project the data on the space generated by the first r eigenvectors $\mathbf{u}^1, \mathbf{u}^2, \dots, \mathbf{u}^r$.

To represent \mathbf{x}^k on this basis we compute the vector $\mathbf{y}^k \in \mathbb{R}^r$ thus achieving dimensionality reduction.

$$\mathbf{y}^k = \begin{pmatrix} (\mathbf{u}^1)^T \mathbf{x}^k \\ (\mathbf{u}^2)^T \mathbf{x}^k \\ \vdots \\ (\mathbf{u}^r)^T \mathbf{x}^k \end{pmatrix} \in \mathbb{R}^r$$

$\mathbf{u}^1, \mathbf{u}^2, \dots, \mathbf{u}^r$ - the principal components

PCA algorithm

1. Compute the mean value $\mu = \frac{1}{p} \sum_{k=1}^p x^k$
2. Subtract mean from data: replace x^k with $x^k - \mu$
/* 3. + 4. can be skipped if data have the same scale
3. Compute $\sigma_i^2 = \frac{1}{p} \sum_{k=1}^p (x_i^k)^2, i = 1, \dots, n$
4. Replace each x_i^k with $\frac{1}{\sigma_i} x_i^k$
5. Compute new $X = \begin{bmatrix} x^1 & x^2 & \dots & x^p \end{bmatrix}$
6. Compute the covariance matrix $C = \frac{1}{p} X * X^T$
7. Compute the eigenvalues and eigenvectors

(orthonormal!) of the covariance matrix

8. Sort the eigenvalues (and eigenvectors) in decreasing order

9. Choose first r eigenvectors $\mathbf{u}^1, \mathbf{u}^2, \dots, \mathbf{u}^r$

$$\mathbf{U} = \begin{bmatrix} \mathbf{u}^1 & \mathbf{u}^2 & \dots & \mathbf{u}^r \end{bmatrix} \in \mathbb{R}^{n \times r}$$

10. Project data on principal components

$$\mathbf{Y} = \mathbf{U}^T * \mathbf{X} \in \mathbb{R}^{r \times p}, \mathbf{Y} = \begin{bmatrix} \mathbf{y}^1 & \mathbf{y}^2 & \dots & \mathbf{y}^p \end{bmatrix}$$

Andrew Ng - *Principal Component Analysis*

<http://cs229.stanford.edu/notes/cs229-notes10.pdf>

Lindsay I. Smith - *A tutorial on principal components analysis*, 2002.

<https://ourarchive.otago.ac.nz/bitstream/handle/10523/7534/OUCS-2002-12.pdf?sequence=1>

Jonathon Shlens - *A tutorial on principal component analysis*, *arXiv preprint arXiv:1404.1100* (2014).

<https://arxiv.org/abs/1404.1100>

Face Recognition – Eigenfaces

(Sorry, the code is “narrative” ☺)

```
cale_att='./ATTfaces/';  
loc_images=dir(strcat(cale_att,'*.pgm'));  
  
X=[];  
i_X=1;  
X_test=[];  
i_X_t=1;  
Nume=zeros(length(loc_images),1);  
Num_poza=zeros(length(loc_images),1);  
for i=1:length(loc_images)  
    numef=strcat(cale_att,loc_images(i).name);  
    I=imread(numef);  
    [m, n]=size(I);
```

```

% label processing
nume=loc_images(i).name;
ind=strfind(nume, '.');
nume=nume(2:ind(1)-1);
if(str2double(nume) < 100)
    label=str2double(nume(1));
    nr_poza=str2double(nume(2:end));
elseif str2double(nume(2:3)) == 10 && length(nume) == 3
    label=str2double(nume(1));
    nr_poza=str2double(nume(2:end));
else
    label=str2double(nume(1:2));
    nr_poza=str2double(nume(3:end));
end;
Nume(i)=label;
Num_poza(i)=nr_poza;

```

```
%linearize images
vI=double(I(:));
% build data matrix test/training
if nr_poza == 1 || nr_poza == 2 || nr_poza == 3
    X_test=[X_test, vI];
    Label_test(i_X_t)=label;
    i_X_t=i_X_t+1;
else
    X=[X, vI];
    Label(i_X)=label;
    i_X=i_X+1;
end;
end;
```

```

% compute average image
ave_im=mean(X,2);
% transform data in data with zero mean
X=X-repmat(ave_im,1,size(X,2));
disp(size(X));
[u,s,ve]=svd(X);
%compute the covariance matrix
disp(' Matricea de covarianta ');
tic;
covx=cov(X');
% compute eigenvalues and eigenvectors
disp(' Calcul valori si vectori proprii');
[V,D] = eig(covx);
saveV=V;
eigval = diag(D);

```



```

% sort eigenvalues (rearrange eigenvectors)
eigval = eigval(end:-1:1);
V = fliplr(V);
[eigval1,ind]=sort(-diag(D));
eigval1=-eigval1;
V1=saveV(:,ind);
disp('Elapsed time '); disp(toc);
% https://en.wikipedia.org/wiki/Eigenface
% compute the number of principal components
sum(eigenvalues_pc)=95%*sum(all eigenvalues)
eigsum = sum(eigval);
csum = 0;
for i = 1:length(eigval)
    csum = csum + eigval(i);
    tv = csum/eigsum;
    if tv > 0.95

```

```
    k95 = i;  
    break;  
end;  
end  
disp(k95);
```

```
% https://blog.cordiner.net/2010/12/02/eigenfaces-face-recognition-matlab/
```

```
num_eigenfaces=10;  
image_dims=[112,92];  
figure;  
for n = 1:num_eigenfaces  
    subplot(2, ceil(num_eigenfaces/2), n);  
    eig_face = reshape(V(:,n), image_dims);  
    imshow(eig_face,[ ]);
```

end

% fast eigenfaces computation

tic;

xtx=X'*X/size(X,2);

[Vf,Df]=eig(xtx);

eigval_f=diag(Df);

[eigval_f,ind]=sort(-eigval_f);

eigval_f=-eigval_f;

Vf=Vf(:,ind);

Eigen_faces=X*Vf;

disp('Elapsed time '); disp(toc);

```
%https://en.wikipedia.org/wiki/Eigenface  
% compute the number of principal components  
sum(eigenvalues_pc)=95%*sum(all eigenvalues)
```

```
eigsum = sum(eigval_f);  
csum = 0;  
for i = 1:length(eigval)  
    csum = csum + eigval_f(i);  
    tv = csum/eigsum;  
    if tv > 0.95  
        k95_f = i;  
        break;  
    end;  
end  
disp(k95_f);
```

```
%figure
```

```
num_eigenfaces=10;
```

```
image_dims=[112,92];
```

```
figure;
```

```
for n = 1:num_eigenfaces
```

```
    subplot(2, ceil(num_eigenfaces/2), n);
```

```
    eig_face = reshape(Eigen_faces(:,n), image_dims);
```

```
    imshow(eig_face,[ ]);
```

```
end
```

%Face recognition with 1NN

k95=100;

Eigenfete=V(1:k95,:);

%Eigenfete=Eigen_faces(:,1:k95_f)';

%Project data on eigenfaces

features_training=Eigenfete*X;

features_test=Eigenfete*(X_test-repmat(ave_im,1,size(X_test,2)));

features_training=features_training';

features_test=features_test';

knn_model=fitcknn(features_training,Label);

predicted_labels=predict(knn_model, features_test);

disp('Recognition Error');

disp(sum(predicted_labels~=Label_test)/length(Label_test)*100);

Web search – Google's PageRank

- informații nestructurate, 'self-organized'
- fără standarde, recenzenti, 'paznici' ai conținutului
- informație volatilă, eterogenă și foarte diversă (conținut, structură, format, limbi, alfabet, scop)
- 1.872.963.735 – pagini web active
- $\approx 45 \cdot 10^9$ – pagini web indexate (probabil mai multe) ???

Căutarea

- *roboți web* – parcurg permanent graful Web pentru informații
- *depozițul de pagini* - memorează temporar paginile web (paginile foarte populare sunt stocate perioade mai lungi)
- *modulul de indexare* – se extrage din paginile depozitate informația esențială (cuvinte-cheie, fraze-cheie, descriptori) → se obține o versiune comprimată a paginii; în funcție de popularitatea paginii aceasta este fie ștearsă, fie păstrată în depozitul de pagini

- informațiile esențiale memorate despre orice pagină Web
 - *index de conținut* - cuvinte sau fraze cheie, titlu (se construiește inverted file - similar indexului unei cărți)
 - *index de structură hyperlink*
 - *index de informații speciale* – imagini, pdf

- *modul de interogare* – transformă cererea din limbaj natural într-un limbaj pentru motorul de căutare (de obicei numere); se consultă informațiile esențiale memorate; se obțin paginile relevante cererii (care conțin termeni din interogare).
- *modulul de ‘ranking’(evaluare a relevanței)* – prelucrează paginile relevante furnizate de modulul de interogare și le ordonează după anumite criterii (funcție de motorul de căutare); acest modul trebuie să discearnă care pagină Web răspunde cel mai bine interogației și le ordonează;

PageRank este sistemul de evaluare a paginilor web din punctul de vedere al importanței, folosit de motorul de căutare Google. Nota furnizată de PageRank este unul din factorii folosiți pentru a ordona paginile web atunci cand se face o căutare.

“**PageRank** is an important factor. It is **one out of 200** signals but still it is an **important** one for a large number of queries.”
(Matt Cutts – head of Google’s Webspam team)

Google dă cel puțin 2 ‘note’/’scoruri’ fiecărei pagini Web:

- a. popularitate
- b. conținut – modul de calcul este secret – termenii din interogare sunt în titlu sau în interiorul paginii, de câte ori termenii apar în pagină, cât de aproape unii de alții sunt termenii din interogațiile cu cuvinte multiple, cum apar termenii de interogare (bold, subtitluri...), conținutul paginilor vecine

Nota/scorul finală cu care se face ordonarea este o medie ponderată între notele de mai sus, componenta principală fiind popularitatea.

PageRank (Google) – S. Brin, L. Page (popularitate)

Sergey Brin, Lawrence Page, R. Motwami, Terry Winograd – ‘*The PageRank citation ranking: bringing order to the Web*’ – Technical Report 1999-0120, Computer Science Department, Stanford University, 1999

Ideea: *O pagină este importantă dacă este referită (citată) de alte pagini importante.*

Importanța unei pagini (PageRank-ul ei) este calculată sumând PageRank-urile tuturor paginilor care citează pagina respectivă. Dacă o pagină importantă citează mai multe pagini, atunci importanța ei se împarte egal paginilor pe care le citează.

Graful Web – cu n noduri.

$$I_{ij} = \begin{cases} 1 & \text{dacă pagina } j \text{ citează pagina } i \\ 0 & \text{altfel} \end{cases}$$

c_j = numărul de citări ale paginii j (outlink-urile paginii j)

r_i = PageRank-ul paginii i

$$r_i = (1 - d) + d \sum_{j=1}^n \frac{I_{ij}}{c_j} r_j, \quad d > 0, \quad d = ? \quad \mathbf{0.85} \quad (1)$$

Constanta d asigură faptul că fiecare pagina Web are un Pagerank de cel puțin $(1-d)$.

Relația (1) poate fi scrisă matricial astfel:

$$\mathbf{r} = (1 - d)\mathbf{e} + d \mathbf{e}^T \mathbf{I} \mathbf{C}^{-1} \mathbf{r} \quad (2)$$

$$\mathbf{e} = \begin{pmatrix} 1 \\ 1 \\ \vdots \\ 1 \end{pmatrix} \in \mathbb{R}^n, \quad \mathbf{I} = (\mathbf{I}_{ij}) \in \mathbb{R}^{n \times n}, \quad \mathbf{C} = \text{diag}(c_1, c_2, \dots, c_n) \in \mathbb{R}^{n \times n}$$

Deoarece suma tuturor Pagerank-urilor este $\mathbf{1}$ putem scrie:

$$\sum_{i=1}^n r_i = \mathbf{e}^T \mathbf{r} = \mathbf{1}.$$

$$\mathbf{A} = \frac{1-d}{n} \mathbf{e} \mathbf{e}^T + d \mathbf{I} \mathbf{C}^{-1} \in \mathbb{R}^{n \times n}$$

Relația (2) se poate scrie:

$$\mathbf{r} = A\mathbf{r}$$

Această relație ne spune că matricea A (matrice de dimensiuni foarte mari și rară) are valoarea proprie 1 iar vectorul de Pagerank-uri este vectorul propriu asociat. Se poate folosi metoda puterii pentru determinarea lui \mathbf{r} .