

Separarea surselor

Motrescu Radu

Contents

1	Introducere in separarea surselor //TODO	5
2	PCA	6
3	Kernel PCA	6
3.1	Libraria Accord.NET	6
3.2	Ce este PCA?	7
3.3	Jurnalul dezvoltarii aplicatiei	12
3.4	Aplicatia AccordPCA	12
3.5	Aplicatia PlotPCA	13
3.5.1	Setul 3-cloud-points	13
3.5.2	Setul two-moon si abordarea Kernel PCA	15
3.5.3	Cros-validare 10-fold	17
3.6	Aplicatia Eigenfaces	19
3.6.1	Proiectarea unei imagini noi	20
3.6.2	Clasificarea unei imagini noi	21
4	ICA //TODO	22
4.1	Ce este ICA?	22
4.2	Definitia riguroasa ICA	25

4.3	Limitările ICA	26
4.4	Independenta	27
4.4.1	Definitie	27
4.4.2	Corelatia si independenta	28
4.4.3	Variable nongausiene	28
4.5	Principiile estimarii ICA	29
4.5.1	Nongaussianitate inseamna independenta	29
4.6	Masuratori ale nongaussianitatii	31
4.6.1	Kurtosis / Aplatizarea	31
4.6.2	Negentropia	33
4.6.3	Aproximarea negentropiei	34
4.7	Abordarea Maximum Likelihood	36
4.8	Metoda Gradient Descent	36
4.9	Metodologia si algoritmul	38
4.9.1	Jupyter Notebook	38
4.9.2	Semnalele audio testate	38
4.9.3	Pre-procesarea datelor	38
4.9.4	Aplicarea algoritmului Gradient Descent	39
4.9.5	Conditia de oprire	40
4.9.6	Rezultate	40

5	Concluzii //TODO	41
6	Bibliografie //TODO	41

1 Introducere in separarea surselor //TODO

Aceasta lucrare are scopul de a informa cititorul asupra unor metoda de analizare a componentelor unor date, fie acestea date economice sau imagini, pentru a putea fi interpretate mai usor, cat si metode de separare a componentelor unui set de date.

Algoritmi abordati sunt **PCA (Principal Component Analysis)**, **KPCA (Kernel Principal Component Analysis)** si **ICA (Independent Component Analysis)**.

O scurta descriere a cazurilor de aplicare a algoritmilor mentionati:

- PCA - Este folosit pentru a extrage informatiile cele mai relevante dintr-un set de date, adeseori cu dimensionalitate mare, cu scopul atat de a face interpretarea lor mai usoara, cat si pentru prelucrarea lor mai rapida.
- KPCA - Este o extensie a PCA, si acopera slabiciunile PCA, spre exemplu in seturi de date in care distributia valorilor nu poate fi analizata la fel de bine de PCA, si in cazuri de clustering.
- ICA - Aceasta abordare de analizare a componentelor are avantajul de le putea si separa, folosind atat unelte matematica cat si statistice, si este adeseori folosita in cazuri de separare a semnalelor de diferite tipuri, precum audio si video.

2 PCA

3 Kernel PCA

3.1 Libraria Accord.NET

Libraria Accord.Net contine clase pentru:

- Calcul stiintific: matematica, statistica si machine learning
- Procesare de imagini si semnale: imagini, semnale audio si recunoastere si urmarire faciala in timp real
- Biblioteci suport pentru controale specifice: histograme, scatter-plots, controale pentru fiecare clasa de procesare imagini si semnale

In contextul cerintelor pe care vrem noi sa le indeplinim, vom folosi clase dedicate matematicii, statisticii, si ale unor controale de afisare a datelor.

3.2 Ce este PCA?

PCA, sau Principal Component Analysis, sau pe romaneste, Analiza Componentelor Principale este o unealta matematica aplicata din algebra liniara.

Este o metoda non-parametrica (care nu depinde de statistici) de a extrage informatiile relevante dintr-un set de date complex sau confuz.

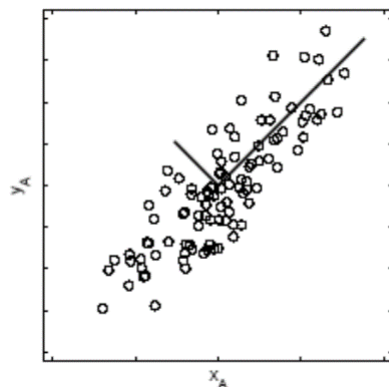
3.2.0.1 Un scurt exemplu: Sa presupunem ca am extras informatii pentru 100 de parametri pentru un student: inaltime, varsta, greutate, nota obtinuta la un test, culoarea parului etc. Vrem sa gasim cele mai importante caracteristici care definesc studentul. Cum facem asta? Folosim PCA pentru a selecta numai cele mai importante caracteristici.

- Ce parametri dorim sa indepartam:
 - Constanti: numarul de capete, care este 1 pentru toti studentii
 - Aproape constanti: grosimea firului de par: 0.003, 0.002, 0.0005 etc.
 - Care depind de alti parametri
- Ce parametri dorim sa pastram:
 - Care nu depind de alti parametri: culoarea ochilor
 - Care se schimba mult, au variatie mare: notele

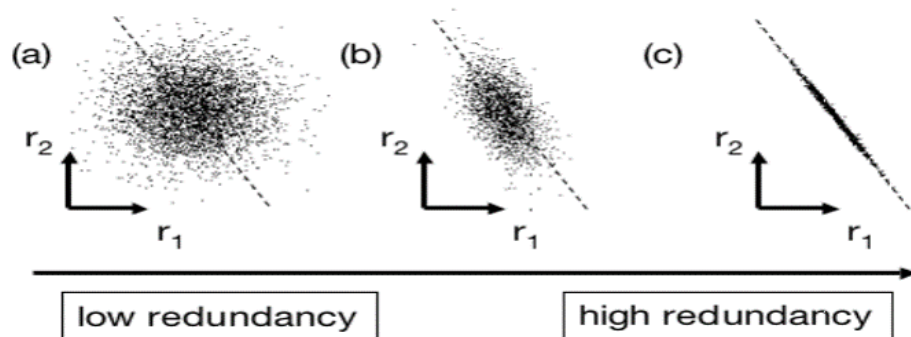
A putea elimina parametrii irelevanti si a-i pastra pe cei relevanti este usor pentru un om, el putand vedea clar care parametri nu exprima informatii relevante despre subiectul respectiv, dar cum putem face calculatorul sa isi dea seama de acesti parametri? Folosind matematica, desigur!

Dorim sa minimizam "sunetul de fundal" si redundanta datelor si sa maximizam variatia dintre parametri.

Se poate vedea in imaginea de mai jos maximizarea variatiei dintre axele norului de puncte respectiv.



In imaginile de mai jos se pot vedea inregistrările unor informatii. In imaginea (a) si (b) se poate vedea cum informatiile nu sunt corelate, avand redunanta mica spre medie (ex: inaltimea unui student si media lui), iar in imaginea (c) se poate vedea o redunanta mare, insemnand ca ambii parametrii pot fi exprimati unul in functie de celalalt.



3.2.0.2 Variatia Este un mijloc de realizare a variabilitatii datelor dintr-un set de date cu media \bar{X} :

$$\sigma^2 = \frac{\sum_{i=1}^n (X_i - \bar{X})^2}{n - 1} \quad (1)$$

3.2.0.3 Covariatia Reprezinta variabilitatea fiecarei dimensiuni in relatie cu celelalte, si este masurata intre 2 dimensiuni pentru a se putea observa relatia dintre cele 2, spre exemplu numarul de ore studiate si nota obtinuta la examen.

$$\text{var}(X) = \frac{\sum_{i=1}^n (X_i - \bar{X})(X_i - \bar{X})}{n - 1} \quad (2)$$

$$\text{cov}(X, Y) = \frac{\sum_{i=1}^n (X_i - \bar{X})(Y_i - \bar{Y})}{n - 1} \quad (3)$$

3.2.0.4 Matricea de covariatie Consideram setul de date din care extragem valoarea medie (zero-mean data) si avand setul de vectori $\{x_1, x_2, \dots, x_m\}$ care reprezinta liniile unei matrici $X_{m,n}$.

Fiecare linie a matricii reprezinta toate masuratorile unui anumit parametru, iar fiecare coloana reprezinta masuratorile care s-au intamplat la un moment dat.

De aici ajungem la definitia matricei de covariatie:

$$S_x \equiv \frac{1}{n - 1} X X^T \text{ where } X = \begin{bmatrix} x_1 \\ \vdots \\ x_m \end{bmatrix} \quad (4)$$

S_x este o matrice simetrica $m \times m$, termenii de pe diagonala reprezentand variatia din acelasi parametru, iar termenii care nu sunt de pe diagonala reprezinta covariatia dintre parametri diferiti. Calculand S_x , cuantificam corelatia dintre toate posibilele perechi de masuratori. Observand elementele

din matrice, o covariatie mare reprezinta un caz de redundanta mare, iar o covariatie egala cu 0 reprezinta date complet necorelate.

$$C = \begin{bmatrix} cov(X, X) & cov(X, Y) & cov(X, Z) \\ cov(Y, X) & cov(Y, Y) & cov(Y, Z) \\ cov(Z, X) & cov(Z, Y) & cov(Z, Z) \end{bmatrix} \quad (5)$$

3.2.0.5 Valorile si vectorii proprii Urmatorul pas in calcularea PCA este aflarea valorilor si a vectorilor proprii ale matricii de covariatie. Extragand aceste informatii, ele ne vor arata componentele principale ale setului de date: vectorul propriu cu cea mai mare valoare proprie este componenta principala a setului de date. Se obisnuieste sortarea vectorilor proprii in functie de valoarea proprie pentru a determina ordinea de semnificativitate.

Vectorii si valorile proprii reies din probleme de urmatoarea forma:

$$A.v = \lambda.v \quad (6)$$

A: matrice $m \times m$

v: vector $m \times 1$ nenul

λ : constanta

Pentru orice valoare a lui λ pentru care ecuatia are solutie se numeste valoarea proprie a lui A, si vectorul **v** care corespunde acestei valori se numeste vectorul propriu a lui A.

3.2.0.6 Pasul final PCA este sa aflam valorile finale ale setului de date. Aici avem optiunea sa ignoram o parte din dimensiunile pe care le avem, deoarece ele pot fi nesemnificative, avand valori proprii mici, sau dorim afisarea lor in 1, 2 sau 3 dimensiuni. Dupa stabilirea vectorilor proprii doriti, aflarea datelor finale este simpla, proiectam punctele in spatiul lor:

$$\text{FinalData} = \text{RowFeatureVector} \times \text{RowZeroMeanData}$$

RowFeatureVectore este matricea cu vectorii proprii transpusi, cu cel mai semnificati vector propriu pe prima linie.

RowZeroMeanData este matricea care contine setul de date initial din care s-a scazut valoarea medie.

Matricea rezultata **FinalData** va contine setul de date dupa aplicarea algoritmului PCA.

3.3 Jurnalul dezvoltarii aplicatiei

3.4 Aplicatia AccordPCA

Primii pasi pe care i-am facut in cadrul acestui proiect a fost implementarea algoritmului PCA. Acest lucru a fost realizat usor, folosind libraria Accord, mai specific, pentru a realiza diferitele operatii cu matrici.

```
public void Compute()
{
    mean = initialData.Mean(0);
    dataAdjusted = initialData.Subtract(mean, 0);
    covarianceMatrix = dataAdjusted.Covariance();

    evd = new EigenvalueDecomposition(covarianceMatrix);

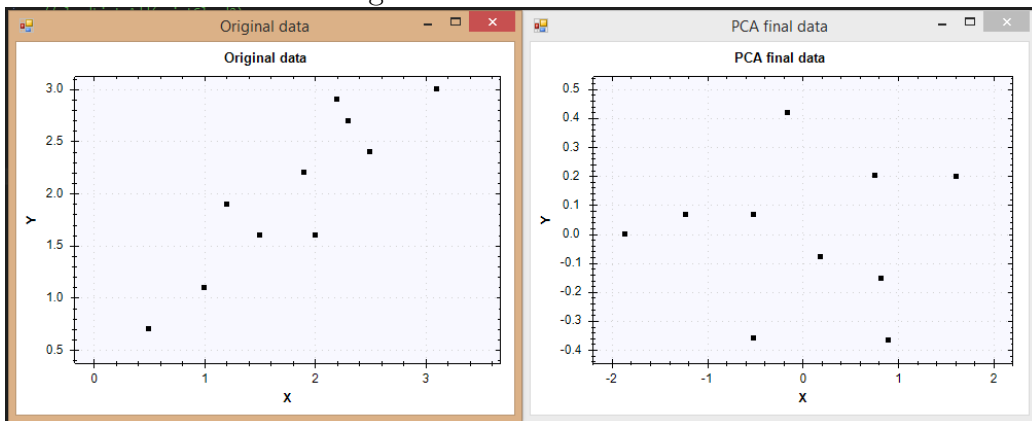
    eigenvalues = evd.RealEigenvalues;
    eigenvectors = evd.EigenVectors;

    eigenvectors = Matrix.Sort(eigenvalues, eigenvectors, new GeneralComparer(ComparerDirection.Descending, true));

    finalData = dataAdjusted.Dot(eigenvectors);
    // data2D = finalData.GetColumns(0, 1);
}
```

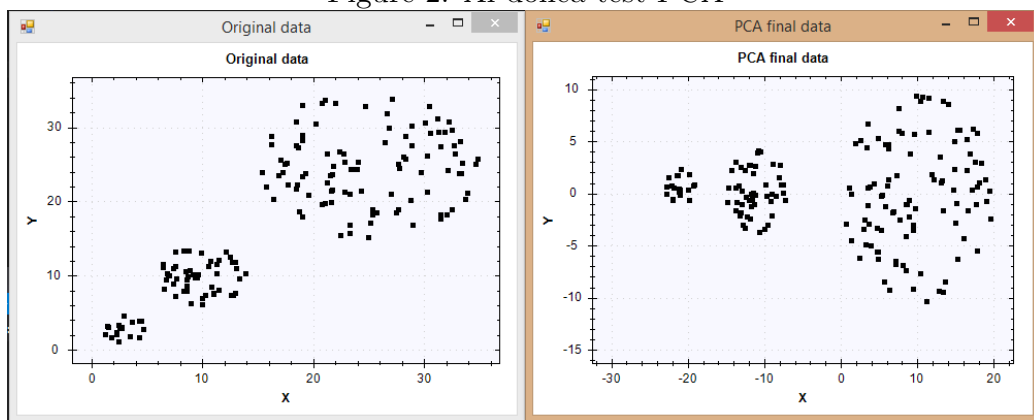
Pentru a testa acest algoritm, am introdus prima data un set de date deja testat de catre altcineva, unde se poate observa clar rezultatul algoritmului. (referinta catre <http://dai.fmph.uniba.sk/courses/ml/sl/PCA.pdf>)

Figure 1: Primul test PCA



Urmatorul lucru pe care l-am dorit sa il facem a fost sa vedem cum mai multe noruri de puncte, de diferite dimensiuni, se modifica dupa aplicarea algoritmului. Pentru acest lucru, am creat o clasa **PointCloud** care descrie un nor de puncte. In interiorul clasei am generat norul de puncte in jurul unor coordonate (x, y) , sub forma unui disc, in care punctele au distributie normala. Rezultatele au fost urmatoarele:

Figure 2: Al doilea test PCA



La fel ca in primul test, se poate observa aranjarea dupa prima componenta principala a norurilor.

3.5 Aplicatia PlotPCA

In acest moment, am decis sa realizam o aplicatie care are un GUI usor de folosit si unde fiecare nor de puncte este reprezentat in alta culoare, pentru a se vedea mai clar ce se intampla, acest lucru fiind mai important pentru seturile de date in care masuratorile pentru fiecare entitate sunt intrepatruse.

3.5.1 Setul 3-cloud-points

Primele teste pe care le-am facut, au fost pe 3 nori de puncte generati aleator, la fel ca in exemplul de mai sus, si se poate vedea mai clar separarea

norilor atat pe in spatiu 2D, cat si pe axa primei componente, cea cu cea mai mare variatie. De asemenea, atunci cand norii se suprapun putin de-a lungul axei Y, algoritmul PCA reuseste sa separe norii complet de-a lungul axei primei componente. Pe de alta parte, daca norii se suprapun mai mult, atunci se poate observa ca algoritmul PCA nu mai poate separa norii la fel de bine ca si pana acum.

Figure 3: Nori separati total

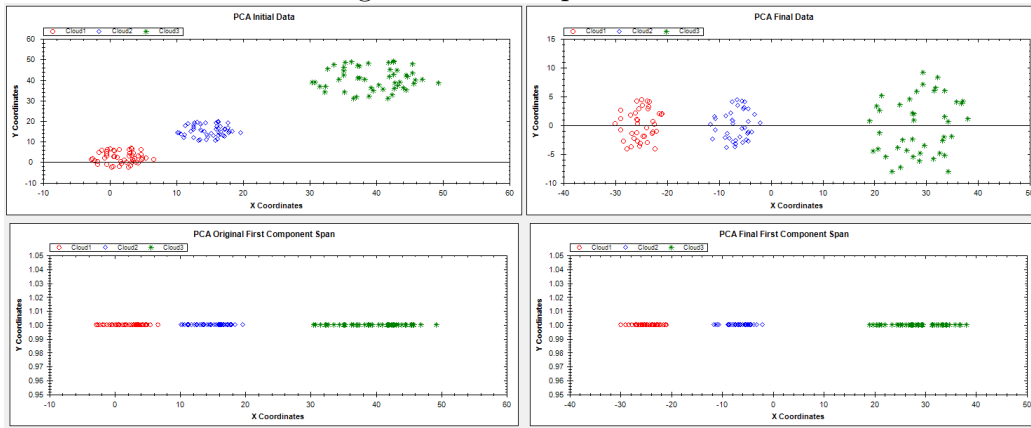


Figure 4: Nori interesectati partial

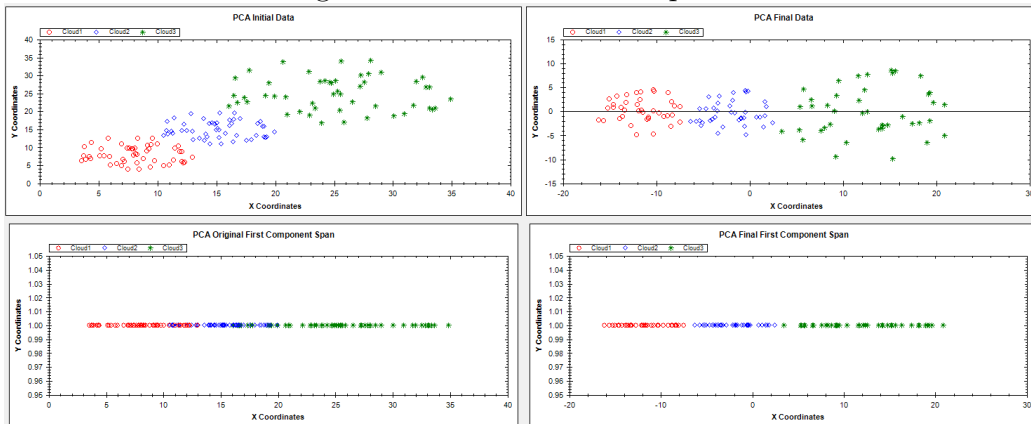
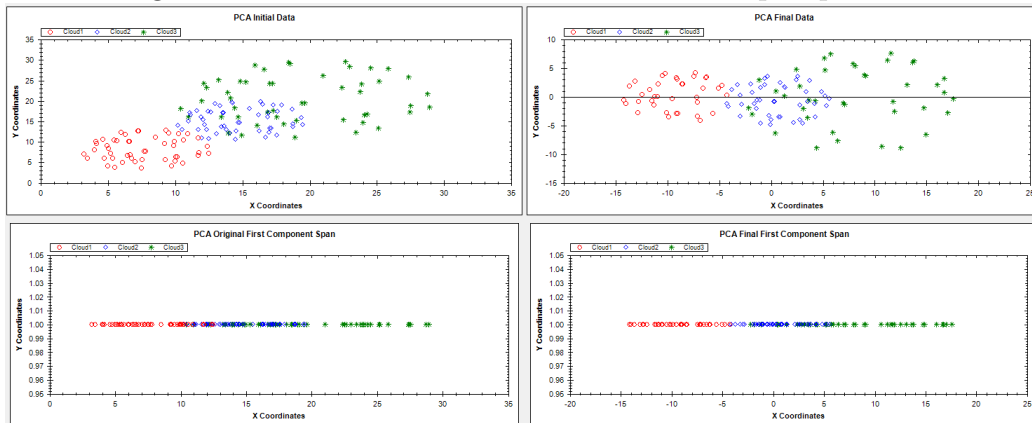


Figure 5: Norul albastru si verde intersectati aproape total



3.5.2 Setul two-moon si abordarea Kernel PCA

Pentru a aborda problema seturilor de date liniar inseparabile, nu putem folosi abordarea PCA simpla, rezultatul nu va fi cel dorit, dupa cum se poate vedea aplicand algoritmul PCA pentru setul two-moon: cele doua seturi de puncte se amesteca de-a lungul primei componente principale, astfel o clasificare a unui nou punct proiectat in spatiul respectiv va fi dificila, deoarece nu putem seta un punct de delimitare intre cei doi nori, acestia fiind intrepunsi.

Abordarea Kernel PCA rezolva acest lucru prin proiectarea datelor intr-un spatiu dimensional mai mare, unde ele vor deveni liniar separabile, iar acest lucru se face cu o "functie kernel".

Pentru a implementa Kernel PCA, vom avea in considerare urmatoarele 2 lucruri:

1. Calcularea matricii kernel

Pentru fiecare pereche de puncte se calculeaza :

$$k(x_i, x_j) = \exp(-\gamma \|x_i - x_j\|_2^2) \quad (7)$$

Daca avem un set de date cu 100 de mostre, matricea kernel va fi o matrice simetrica 100×100 .

Alegerea valorii lui γ este foarte importanta, in functie de ea rezultatul va fi cel dorit sau nu.

2. Calcularea vectorilor si a valorilor proprii

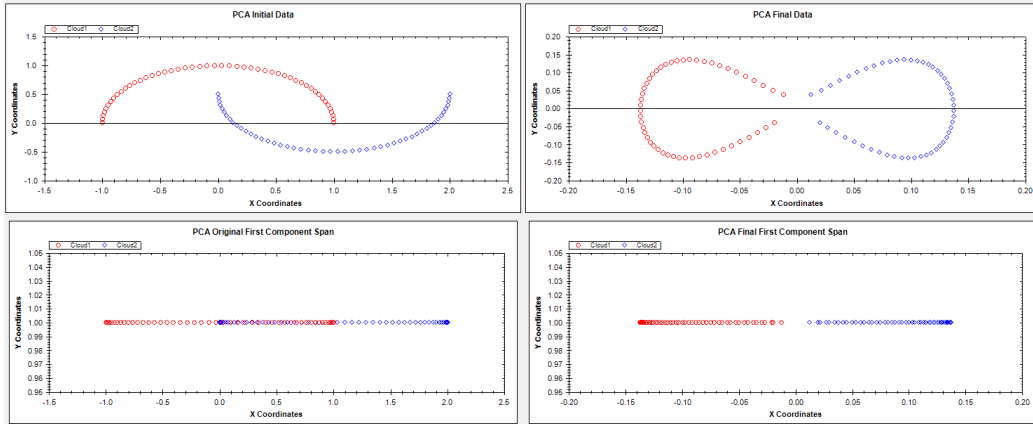
Deoarece nu putem garanta ca matricea e centrata, vom aplica urmatoarea formula pentru a o centra:

$$K' = K - 1_N K - K 1_N + 1_N K 1_N \quad (8)$$

unde 1_N este o matrice $N \times N$ cu toate valorile egale cu $\frac{1}{N}$. Acum, putem afla vectorii si valorile proprii pentru matricea centrata, iar vectorii proprii vor reprezenta setul de date proiectat pe componentele principale respective.

In imaginile urmatoare se poate observa cum algoritmul PCA si algoritmul KPCA afecteaza setul de date two-moon, rezultatul obtinut prin folosirea KPCA este mult mai bun.

Figure 6: Setul two-moon proiectat cu ajutorul algoritmul KPCA, $\gamma = 15$



3.5.3 Cros-validare 10-fold

Pentru a ne asigura ca algoritmul se comporta corect si in alte situatii, am decis sa facem o cros-validare 10-fold (10-fold crossvalidation). Acest lucru implica urmatoarele:

- avem 100 de puncte, distribuite in setul de date two-moon

- din acestea, vom alege pe rand, cate 10 puncte, diferite de fiecare data, si vom retine norul din care face parte fiecare

- cu celelalte 90 de puncte vom "antrena" algoritmul pentru a putea proiecta celelalte 10 puncte pe componentele principale respective, in cazul nostru, dorim sa facem proiectarea pe prima componenta principala

- avand cele 90 de puncte proiectate in spatiul nou, putem determina din nou granitele celor 2 nori de puncte, iar mijlocul acestor granite va fi punctul nostru de separare a norilor

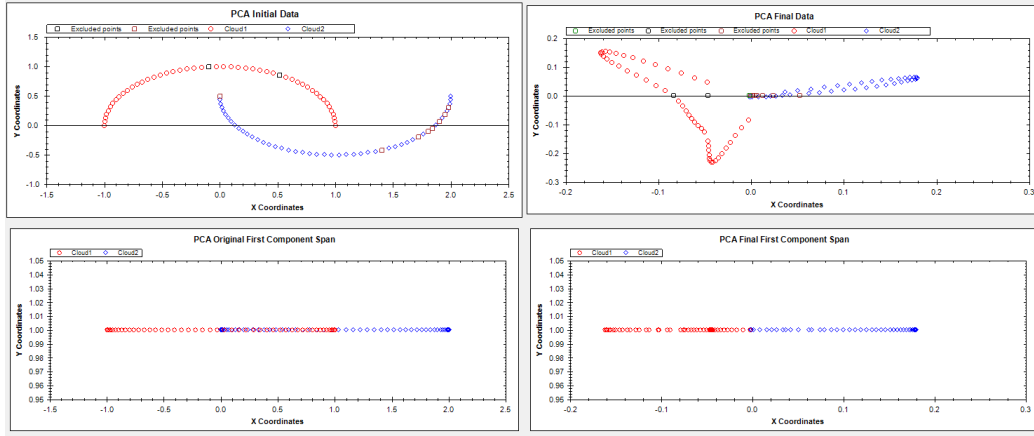
- cu punctul de separare gasit, cele 10 puncte proiectate pe prima componenta principala pot fi categorizate: daca se afla in stanga punctului de separare el va fi in norul albastru de puncte, iar daca este la dreapta punctul de separare, va apartine norului rosu

- stiind din ce nori au provenit punctele si in ce nori au fost proiectati, putem verifica eroarea algoritmului pentru cele 10 puncte:
$$\frac{\text{nr. puncte proiectate gresit}}{\text{nr. puncte totale}}$$

- repetam acest proces pana cand toate punctele au fost proiectate, iar eroarea de final este cea cautata

In cazul testelor noastre, eroarea a variat intre 1% si 4%, fiind un rezultat foarte bun.

Figure 7: Exemplu de validare 10-fold



3.6 Aplicatia Eigenfaces

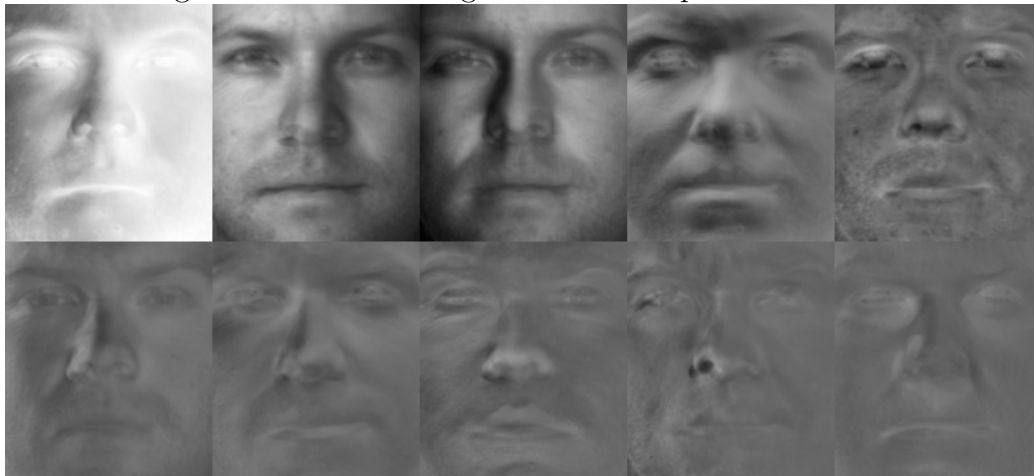
Dupa ce am realizat testele KPCA, urmatorul pas a fost sa crestem dimensiunile cu care lucram. Mai exact, am inceput sa folosim imagini pentru a vedea cum actioneaza algoritmul PCA pe ele.

Imaginile folosite au fost cele din baza de date Yale, care a fost folosita in multe alte cercetari de genul acesta, continand destui subiecti si destule imagini pentru a putea testa orice algoritm.

Pentru a putea procesa aceste imagini, ele trebuiau sa aiba toate aceeasi dimensiune, in cazul nostru 168×192 , iar fiecare imagine va deveni un rand din matricea care va fi prelucrata de algoritmul PCA.

Fetele rezultate din acest algoritm reprezinta toti vectorii proprii care contribuie la spatiul fetelor intr-un mod semnificativ, si sunt sortati dupa valoarea proprie, si dupa cum se poate observa, aceste imagini devin din ce in ce mai nesemnificative:

Figure 8: Primele 10 eigenfaces-uri ale primului subiect



3.6.1 Proiectarea unei imagini noi

Pentru a putea proiecta o imagine noua in spatiul fetelor, tot ce trebuie sa facem este:

$$X = X - avg(X) \quad (9)$$

$$Y = X^T \times E_0 \quad (10)$$

X = este vectorul imaginii pe care o vom proiecta

E_0 = este vectorul propriu asociat valorii proprii maxime

Y = este imaginea proiectata.

3.6.2 Clasificarea unei imagini noi

Pentru a vedea daca o imagine nou proiectata apartine sau nu setului initial de imagini, avem nevoie mai intai de imaginea proiectata, si de o granita, la fel ca la abordarea KPCA.

Modul prin care vom afla daca o imagine apartine setului este prin calcularea distantei de la aceasta imagine la reprezentarea setului de imagini in spatiul respectiv, iar daca valoarea aceasta este mai mica decat granita mentionata mai sus, imaginea va fi clasificata ca va apartine setului initial.

Pentru a afla aceasta distanta, vom efectua urmatoarele operatii:

$$D = W_i - Y \quad (11)$$

$$N = \|D_i\| \quad (12)$$

$$v = \min(N) \quad (13)$$

W = reprezinta fetele setului de date initial proiectate in spatiul fetelor

D = diferenta dintre toate liniile lui W si Y N = norma fiecărei linii din D

v = este distanta dintre imaginea nou proiectata si spatiul fetelor

Pentru a afla granita care delimiteaza imaginile din spatiul fetelor de celelalte, vom folosi aceleasi operatii ca si mai sus, dar aplicate pe toate imaginile setului initial, si vom cauta $v = \text{avg}(N)$, distanta medie dintre toate fetele proiectate si spatiul fetelor. Am ales sa facem acest lucru si sa nu cautam $v = \max(N)$ deoarece se poate intampla ca o imagine din setul initial sa aiba distanta fata de spatiul fetelor disproportional mai mare fata de celelalte, si astfel granita de clasificare ar fi prea mare, incluzand orice imagine pe care incercam sa o proiectam.

Acestea fiind zise, am facut urmatorul test: din cele 65 de imagini ale subiectului 1 din baza de date Yale, primele 44 au fost folosite pentru a "antrena" algoritmul, iar ultimele 21 au fost folosite pentru clasificare. Eroarea rezultata a fost de $\sim 50\%$, imaginile care nu au fost clasificate ca facand parte

in set au fost cele in care subiectul are numai jumatate de fata iluminata, numarul acestor imagini fiind $\sim 50\%$ din ele, deci putem trage concluzia ca algoritmul face ceea ce isi propune.

4 ICA //TODO

4.1 Ce este ICA?

Una dintre problemele din lumea reala este gasirea unei reprezentari potrivite pentru date multivariate, adica date care au multe dimensiuni, si care de multe ori sunt aleatoare. Din cauza problemelor de performanta, orice fel de reprezentare este de obicei interpretata ca o transformare liniara a datelor originale. Printre aceste metode de analiza a transformarilor liniare se numara si PCA, despre care am vorbit pana acum. ICA (Independent Component Analysis) sau Analiza Componentelor Individuale este o abordare noua care are rolul de a reprezentarile liniare ale unor date cu distributie non-gaussiana, astfel incat componentele sa fie cat mai independente din punct de vedere statistic. Acest tip de separare este dorita deoarece captureaza structura esentiala a datelor in multe tipuri de aplicatii, precum extragerea caracteristicilor (feature extraction) sau separarea semnalelor (signal separation).

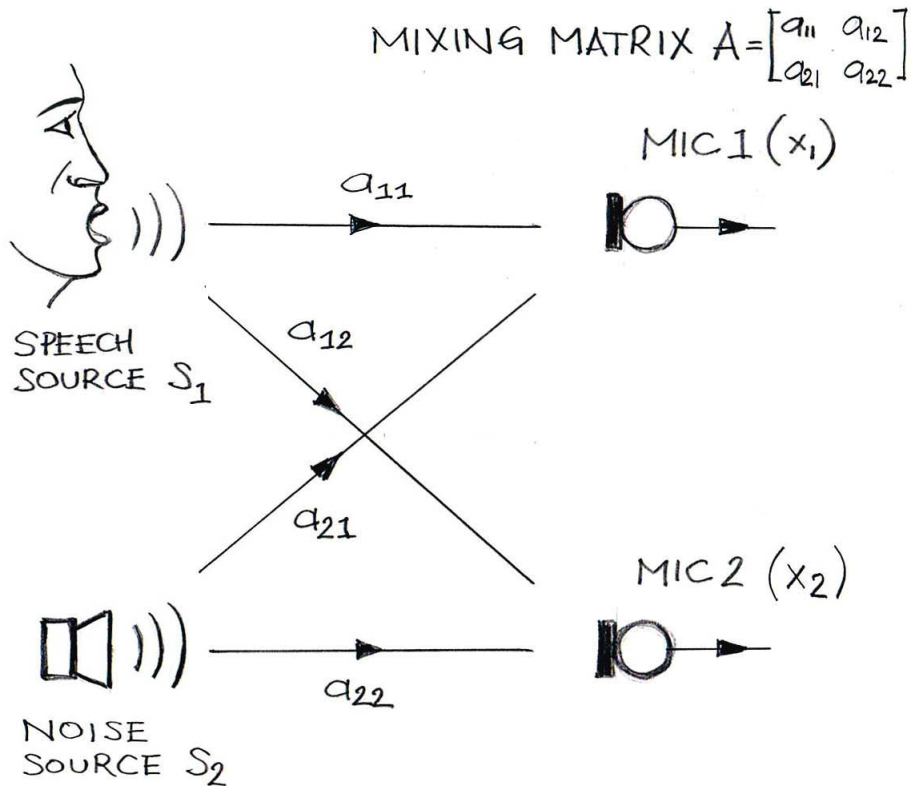
Pentru a putea vizualiza aceasta idee, putem analiza o problema din lumea reala. Sa spunem ca suntem intr-o camera in care 2 oameni vorbesc simultan. Avem 2 microfoane pe care le positionam in locuri diferite in camera, si ele ne dau 2 inregistrari ale semnalelor. Aceste semnale sunt $x_1(t)$ si $x_2(t)$, cu x_1 si x_2 functii de amplitudine in functie de t . Fiecare inregistrare este o suma ponderata a semnalelor emise de cei doi oameni care vorbesc, pe care le notam $s_1(t)$ si $s_2(t)$. Pentru claritate le punem sub forma de ecuatie liniara:

$$x_1(t) = a_{11}s_1 + a_{12}s_2 \tag{14}$$

$$x_2(t) = a_{21}s_1 + a_{22}s_2 \tag{15}$$

unde $a_{11}, a_{12}, a_{21}, a_{22}$ sunt niste parametrii care determina distanta dintre microfoane si vorbitori. Avand aceste date, daca am putea estima cele doua semnale initiale $s_1(t)$ si $s_2(t)$ folosind doar semnalele inregistrate $x_1(t)$ si $x_2(t)$, am reusi sa rezolvam multe probleme din lumea reala. Tipul de problema prezentat mai sus este denumit si *cocktail-party problem*.

In imaginea urmatoare avem o reprezentare vizuala a problemei:



Stiind x_1 si x_2 , daca am cunoaste elementele matricii A , atunci am putea afla s_1 si s_2 prin mijloace normale, fiind o ecuatie liniara. Dar, formuland problema ca mai inainte, informatiile initiale sunt cele 2 semnale inregistrate de la cele 2 microfoane.

O abordare de rezolvare a problemei ar putea fi sa incercam sa ne folosim de proprietatile statistice ale semnalelor $s_1(t)$ si $s_2(t)$ pentru a estima matricea A . Tot ce trebuie sa facem este sa presupunem ca semnalele $s_i(t)$ sunt independente statistic la fiecare moment t .

4.2 Definitia riguroasa ICA

Pornim de la presupunerea ca observam n amestecari liniare x_1, x_2, \dots, x_n a n componente independente:

$$x_j = a_{j1}s_1 + a_{j2}s_2 + \dots + a_{jn}s_n, \text{ pentru toate } j. \quad (16)$$

Se observa ca am renuntat la indicele de timp t , deoarece ICA presupune ca atat fiecare amestecare x_j cat si fiecare sursa s_k sunt variable aleatoare, lucru care o sa fie folositor cand vom vorbi despre distributia datelor. De asemenea, ne vom asigura ca variabilele observate x_j sunt centrate, scazand media esantioanelor.

Este convenient sa ne folosim de notatia vector-matrice, in locul sumelor precum cea precedenta. De acum vom nota \mathbf{x} vectorul al carui elemente sunt amestecarile x_j , si la fel pentru \mathbf{s} , ale carui elemente sunt s_j . De asemenea, vom nota \mathbf{A} matricea de elemente a_{ij} . Folosind aceasta notatie, modelul de amestecare se scrie astfel:

$$\mathbf{x} = \mathbf{A}\mathbf{s}. \quad (17)$$

Modelul prezentat este modelul ICA, si este un model generativ, descris in statistica astfel:

$$P(X|Y = y) \quad (18)$$

unde X este o variabila observabila, iar Y este variabila cautata, ceea ce inseamna ca descrie cum datele observate, x_i , sunt generate de un proces aplicat asupra datelor initiale, in cazul nostru componentele s_i . Componentele independente sunt variabile latente, nu sunt observate direct, ci sunt deduse din alte variabile. De asemenea, presupunem ca matricea de amestecare A nu este cunoscuta. Tot ce observam este vectorul aleator \mathbf{x} , si trebuie sa estimam atat \mathbf{A} cat si \mathbf{s} . Acest lucru trebuie facut sub anumite ipoteze.

Punctul de inceput al ICA este ipoteza ca componentele independente s_i sunt independente statistic, lucru pe care il vom defini mai incolo. Vom presupune deasemenea ca acestea au distributie *nongaussiana*, si ca matricea A de amestecare este patrata. Dupa ce am estimat matricea A , ii putem calcula inversa, sa zicem \mathbf{W} , si obtinem componentele independente astfel:

$$\mathbf{s} = \mathbf{W}\mathbf{x} \quad (19)$$

4.3 Limitările ICA

Având în vedere modelul dat, se pot observa 2 limitări:

1. Nu putem determina *energiile* componentelor independente.

Motivul este că, atât \mathbf{s} cât și \mathbf{A} sunt necunoscute, orice multiplicator scalar al unei dintre surse s_i ar putea fi negat prin împărțirea coloanei corespunzătoare a_i al \mathbf{A} cu acest scalar. Același lucru se întâmplă cu semnul variabilelor, acesta se poate înmulți cu -1 fără a afecta modelul. Din ferici, această limitare nu este semnificativă în majoritatea cazurilor.

2. Nu putem determina ordinea componentelor independente:

Motivul este că, din nou, atât \mathbf{s} cât și \mathbf{A} fiind necunoscute, putem schimba ordinea termenilor din suma precedentă, oricare dintre componente putând fi prima. Formal, având o matrice de permutare \mathbf{P} și inversa ei, modelul poate fi formulat astfel:

$$x = AP^{-1}Ps. \quad (20)$$

Elementele lui \mathbf{Ps} sunt componentele independente originale, dar în ordine diferită, iar matricea \mathbf{AP}^{-1} este doar o nouă matrice de amestecare necunoscută, care poate fi aflată.

4.4 Independenta

4.4.1 Definitie

Ca sa putem intelege conceptul de independenta, consideram 2 variabile scalare aleatoare, y_1 si y_2 . Variabilele y_1 si y_2 sunt independente daca informatia din y_1 nu influnteaza informatia din y_2 si invers. Acest lucru este valabil pentru variabilele precedente s_1, s_2 , dar nu si pentru amestecarile x_1, x_2 .

Independenta poate fi definita prin probabilitati de densitate. O functie de densitate a unei probabilitati este folosita pentru a specifica probabilitatea unei variabile aleatoare de a avea valori dintr-un anumit interval. Probabilitatea este data prin integrarea functiei cu limitele intervalului. Sa notam $p(y_1, y_2)$ functia de densitate cumulata a y_1 si y_2 , si $p_1(y_1)$ functia de densitate marginala a y_1 , care reprezinta functia de densitate a variabilei y_1 :

$$p_1(y_1) = \int p(y_1, y_2) dy_2, \quad (21)$$

si similar pentru y_2 . Acestea fiind spuse, definim y_1 si y_2 ca fiind independente daca:

$$p(y_1, y_2) = p_1(y_1)p_2(y_2) \quad (22)$$

Definitia aceasta poate fi folosita pentru a demonstra una dintre cele mai importante proprietati ale variabilelor aleatoare independente. Folosind functii h_1 si h_2 de mai sus, avem:

$$Eh_1(h_1), h_2(y_2) = Eh_1(y_1)Eh_2(y_2) \quad (23)$$

Putem demonstra astfel:

$$\begin{aligned} E\{h_1(y_1), h_2(y_2)\} &= \iint h_1(y_1)h_2(y_2)p(y_1, y_2)dy_1dy_2 \\ &= \iint h_1(y_1)p_1(y_1)h_2(y_2)p_2(y_2)dy_1dy_2 \\ &= \int h_1(y_1)p_1(y_1)dy_1 \int h_2(y_2)p_2(y_2)dy_2 \\ &= E\{h_1(y_1)\}E\{h_2(y_2)\} \end{aligned} \quad (24)$$

4.4.2 Corelatia si independenta

O forma mai slaba de independenta este necorelatia. Doua variabile independente y_1 si y_2 se numeste necorelate daca covariatia lor este 0:

$$E\{y_1 y_2\} - E\{y_1\}E\{y_2\} = 0 \quad (25)$$

Daca variabilele sunt independente, ele sunt necorelate, rezultand $h_1(y_1) = y_1$ si $h_2(y_2) = y_2$ din ecuatia precedenta. Pe de alta parte, necorelarea nu implica independenta. Sa presupunem ca (y_1, y_2) sunt valori discrete si au o distributie astfel incat perechea are $1/4$ sanse sa ia una dintre valorile urmatoare: $(0, 1), (0, -1), (1, 0), (-1, 0)$. Atunci y_1 si y_2 sunt necorelate, dar:

$$E\{y_1^2, y_2^2\} = 0 \neq \frac{1}{4} = E\{y_1^2\}E\{y_2^2\} \quad (26)$$

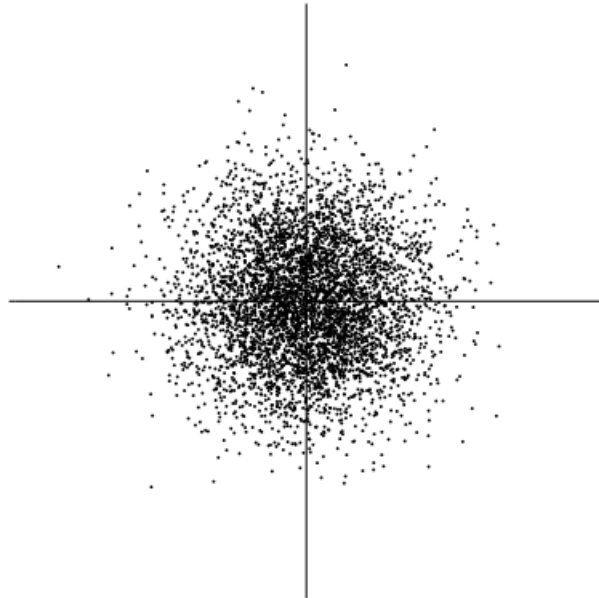
si astfel conditia de mai sus nu este respectata, variabilele nefiind independente.

4.4.3 Variabile nongausiene

Restrictia fundamentala a ICA este ca componentele independente sa fie nongausiene pentru ca separarea sa se intample. Ca sa vedem de ce acest lucru este necesar, sa presupunem ca sursele s_1, s_2 au distributie gaussiana, iar A este o matrice ortonormala. Din asta rezulta ca x_1, x_2 vor fi gaussiene, necorelate si au variatie unitara. Densitatea cumulata a x_1, x_2 este:

$$p(x_1, x_2) = \frac{1}{2\pi} \exp\left(-\frac{x_1^2 + x_2^2}{2}\right) \quad (27)$$

Distributia aceasta este ilustrata mai jos:



Se poate observa cum densitatea este complet simetrica. Asta inseamna ca nu contine informatii cu privire la directiile coloanelor matricii de amestecare A . Din acest motiv, A nu poate fi estimat.

4.5 Principiile estimarii ICA

4.5.1 Nongaussianitate inseamna independenta

Daca luam in calcul cele spuse pana acum, putem deduce ca esenta estimarii unui model ICA sta in nongaussianitate.

Teorema limita centrala, in versiunea clasica, stabileste ca sumele parțiale "normalizate" ale unui sir de variabile aleatoare independente si identic distribuite, cu dispersie finita, tind in distributie catre legea normala (gausiana).

Sa presupunem ca datele din vectorul x este distribuit asemenea modelului ICA prezentat mai devreme, adica ca o amestecare de componente independente, si distributia componentelor independente este identica. Pentru a estima componentele independente, trebuie sa aplicam operatii liniare asupra x_i :

$$y = w^T x = \sum_i w_i x_i, \text{ unde } w \text{ este vectorul pe care il cautam.} \quad (28)$$

Daca w ar fi un rand al inversei lui A , aceasta combinatie liniara ar fi egala cu una dintre componentele independente. Ca sa vedem cum teorema limita centrala este folosita pentru a gasi un w ca sa respecte conditia precedenta, putem face o schimbare de variabile, definind $z = A^T w$. Apoi avem:

$$y = w^T x = w^T A s = z^T s \quad (29)$$

Se observa ca y este o combinatie liniara de s_i , cu ponderile date de z_i . Deoarece suma sumelor celor doua variabile independente este mai gaussiană decat variabilele originale, $z^T s$ este mai gaussiană decat oricare s_i , si devine mai puțin gaussiană atunci când este s_i , în acest caz, doar un element al lui z este diferit de 0.

Cautam w astfel incat $w^T x$ sa fie cat mai nongaussian. Un astfel de vector ar corespunde unui z care are doar o componenta diferita de zero. Asta înseamnă ca $w^T x = z^T s$ este egal cu una dintre componentele independente.

Maximizand nongaussianitatea lui $w^T x$ ne da una dintre componentele independente. Spatiul de cautare pentru optimizarea nongaussianitatii într-un spatiu n -dimensional al vectorilor w are $2n$ maxime locale, 2 pentru fiecare componenta independenta, s_i si $-s_i$. Ca sa gasim mai multe componente independente, trebuie sa gasim toate maximele locale.

4.6 Masuratori ale nongaussianitatii

Pentru a folosi nongaussianitatea in estimarea ICA, avem nevoie de o masura cantitativa a nongaussianitatii unei variabile aleatoare, y . Vom presupune ca y este centrat, adica are media egala cu zero, iar variatia ei este unu.

4.6.1 Kurtosis / Aplatizarea

O masura clasica a nongaussianitatii este aplatizarea, sau al patrulea moment central al unei variabile aleatoare, si este definit astfel:

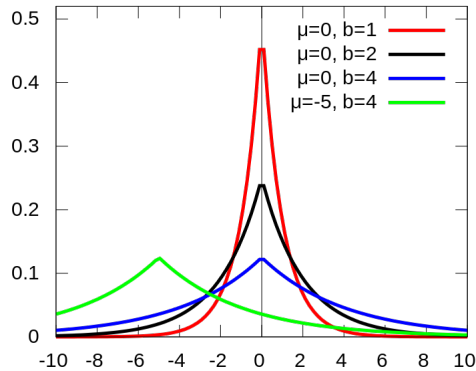
$$kurt(y) = E\{y^4\} - 3(E\{y^2\})^2 \quad (30)$$

Deoarece am stabilit ca variatia este unu, atunci formula se simplifica la:

$$kurt(y) = E\{y^4\} - 3 \quad (31)$$

Acest lucru arata defapt ca aplatizarea reprezinta al patrulea moment, $E\{y^4\}$. Pentru un y gaussian, al patrulea moment este definit ca $3(E\{y^2\})^2$, deci aplatizarea pentru un asemenea y este 0. Pentru majoritatea variabilelor aleatoare nongausiene, aplatizarea este diferita de zero.

Aplatizarea poate fi pozitiva sau negativa, variabilele aleatoare care au aplatizare negativa numinduse subgausiene, iar cele care au aplatizare pozitiva se numest supergausiene. Variabilele aleatoare supergausiene corespund cu distributii de tip Laplace, care au un varf ascutit, si coada alungita:



Variabilele aleatoare subgaussiene corespund cu distributiile uniforme, acestea fiind plate in reprezentare:



Valoarea absoluta a aplatizarii, a fost folosita ca masura a nongaussianitatii in ICA, deoarece este simplu atat teoretic cat si computational. Pentru a calcula aplatizarea, putem folosi pur si simplu al patrulea moment al datelor avute. Analiza teoretica este simplificata datorita proprietatii de liniaritate. Daca avem doua variabile aleatoare independente, atunci:

$$kurt(x_1 + x_2) = kurt(x_1) + kurt(x_2) \quad (32)$$

$$kurt(\alpha x_1) = \alpha^4 kurt(x_1) \quad (33)$$

Pentru a demonstra cum arata spatiul de optimizare pentru aplatizare, si cum se pot afla componentele independente pot fi gasite prin maximizarea sau minimizarea aplatizarii. Avand modelul $x = As$ ca mai inainte, si componentele independente s_1 si s_2 , cu valorile aplatizarii $kurt(s_1)$ si $kurt(s_2)$, ambele diferite de zero, cautam componentele independente prin $y = w^T x$. Facem din nou transformarea $z = A^T w$. Apoi avem:

$$y = w^T x = w^T A s = z^T s = z_1 s_1 + z_2 s_2 \quad (34)$$

Folosind proprietatile aplatizarii mentionate mai devreme avem:

$$kurt(y) = kurt(z_1 s_1) + kurt(z_2 s_2) = z_1^4 kurt(s_1) + z_2^4 kurt(s_2) \quad (35)$$

Mai devreme am spus ca variatia lui y este unu, deoarece s_1 si s_2 au variatia unu. Acest lucru implica faptul ca:

$$E\{y^2\} = z_1^2 + z_2^2 = 1 \quad (36)$$

ceea ce inseamna ca z este restrans la un cerc cu raza unu in spatiul 2D. Problema de optimizare este urmatoarea: care este maximul functiei $|kurt(y)| = |z_1^4 kurt(s_1) + z_2^4 kurt(s_2)|$ pentru acest tip de cerc?

Maximele in acest caz sunt in punctele unde exact unul dintre elementele vectorului z este zero si celalalt este diferit de zero, si deoarece ne aflam intr-un cerc cu raza unu, elementul diferit de zero este egal cu -1 sau 1 . Putem observa ca aceste puncte sunt cele in care y este egal cu unul dintre componentele independente $\pm s_i$, si problema este rezolvata.

In practica am incepe de la un w aleator, am calcula directia in care aplatizarea lui $y = w^T x$ creste cel mai rapid, daca aplatizarea este pozitiva sau scade cel mai rapid, daca aplatizarea este negativa, in functie de esantioanele lui $x(i)$ al vectorului de mixturi x , si folosind o metoda de tip gradient pentru a gasi un nou vector w .

In practica, aceasta metoda are cateva dezavantaje. Aplatizarea este sensibila la valori extreme setului de date, ceea ce ar insemna ca toata aproximarea ar depinde de aceste valori, lucru care cauzeaza rezultate eronate.

4.6.2 Negentropy

Negentropy este a doua metoda foarte importanta pentru masurarea non-gaussianitatii unor variabile aleatoare independente. Din punctul de vedere al teoriei informatiei, negentropy se bazeaza pe cantitatea de entropie diferentiala. Conceptul de entropie este un concept de baza in teoria informatiei. Aceasta reprezinta gradul de informatiei pe care observarea unei variabile o da. Cu cat variabila este mai putin structurata si mai putin previzibila, cu atat mai mare este entropia variabilei. Entropia H a unei variabile aleatoare discrete Y este:

$$H(Y) = - \sum_i P(Y = a_i) \log P(Y = a_i) \quad (37)$$

unde a_i reprezinta valorile posibile alea lui Y . Aceasta definitie poate fi generalizata pentru variabile aleatoare continue, in acest caz, numinduse entropie diferentiala:

$$H(Y) = - \int f(y) \log f(y) dy \quad (38)$$

unde $f(y)$ este functia de densitate a unui vector aleator y .

Un rezultat fundamental al teoriei informatiei este stabilirea faptului ca variabilele gaussiene au cea mai mare entropie, comparativ cu alte variabile cu variatie egala. Asta inseamna ca entropia poate fi folosita ca si masura pentru nongaussianitate, aratand ca variabilele cu distributie gaussiană sunt cele mai puțin structurate dintre toate distributiile. Entropia este mica pentru variabilele ale caror distributii sunt concentrate pe anumite valori, fiind analog unor functii cu varfuri ascutite, de tip Laplace. Pentru a obtine o masura a nongaussianitatii care este zero pentru o variabila aleatoare gaussiană si tot timpul pozitiva, se foloseste o versiune modificata a definitiei entropiei diferentiale, numita negentropie. Negentropia J este definitia astfel:

$$J(y) = H(y_{gauss}) - H(y) \quad (39)$$

unde y_{gauss} este variabila aleatoare gaussiană care are aceasi matrice de covariatie ca si y . Datorita proprietatilor de mai sus, negentropia este pozitiva, si zero atunci cand variabila y are distributie gaussiană.

La fel ca la aplatizare, aproximarea negentropiei are un dezavantaj major, acela ca este greu de calculat din punct de vedere computational, deoarece ar fi nevoie de o estimare a functiei de densitate, lucru care nu poate fi facut precis.

4.6.3 Aproximarea negentropiei

Dupa cum am mentionat mai devreme, aproximarea negentropiei este dificila. In practica, se fac anumite aproximari, iar tipurile de aproximari abordate se va face aici.

$$J(y) \approx \frac{1}{12}E\{y^3\}^2 + \frac{1}{48}kurt(y)^2 \quad (40)$$

Prima abordare este cea bazata pe momentele statistice: Variabila y se presupune ca are media egala cu zero, si variatie unitara. Se poate observa ca aceasta abordare sufera de aceleasi limitari precum aproximarea folosind aplatizarea. Pentru a evita acest lucru, s-au dezvoltat aproximari bazate pe

principiul entropiei maxime, cu forma:

$$J(y) \approx \sum_{i=1}^p k_i [E\{G_i(y)\}] - E\{G_i(v)\}]^2 \quad (41)$$

unde k_i reprezinta constante pozitive, si v este o variabila gaussiană, iar atat v cat si y au media zero si variatie unitara, functiile G_i fiind niste functii nonpatratice. In cazul in care folosim o singura functie nonpatratice, aproximarea devine:

$$J(y) \propto [E\{G(y)\}] - E\{G(v)\}]^2 \quad (42)$$

Alegand G bine, putem obtine aproximari mult mai bune decat cele de mai devreme, spre exemplu, daca alegem G care nu creste repede, precum:

$$G_1(u) = \frac{1}{a_1} \log \cosh a_1 u, G_2(u) = -\exp\left(\frac{-u^2}{2}\right) \quad (43)$$

unde $1 \leq a_1 \leq 2$ este o constanta. Cu toate ca nu am folosit aceste tipuri de aproximari, ele sunt un compromis bun intre metodele de masurare a nongaussianitatii abordate mai devreme, aplatizarea si negentropia.

4.7 Abordarea Maximum Likelihood

O abordare foarte populara pentru rezolvarea modelului ICA este *maximum likelihood estimation*. Aceasta metoda estimeaza valorile parametrului unui model. Acesti parametri sunt gasiti astfel inca estimarea lor sa mareasca sansa ca modelul sa produca datele observate. Pentru cazul particular al modelului nostru ICA, avand $W = (w_1, \dots, w_n)^T$ ca A^{-1} , *log-likelihood-ul* ia forma:

$$L = \sum_{t=1}^T \sum_{i=1}^n \log f_i(w_i^T x(t)) + T \log |\det W| \quad (44)$$

unde f_i este functia de densitate a lui s_i , aici presupunand ca o stim, iar $x(t)$ reprezinta esantioanele amestecarilor x_i . Termenul $\log |\det W|$ vine din faptul ca pentru un vector aleator x cu densitatea p_x , si pentru orice matrice W :

$$y = Wx, \text{ data de } p_x(Wx) |\det W|. \quad (45)$$

4.8 Metoda Gradient Descent

Metoda gradient descent este una dintre cele mai folosite metode pentru aproximarea maximului sau minimului pentru o functie complexa. Acesta functioneaza facand pasi in directia gradientului functiei in punctul curent, astfel se ajunge la minimul functiei. Daca in schimb, pasii sunt in directia pozitiva gradientului, se ajunge la maximul functiei. Gradientul unei functii reprezinta derivata partiala a functiei pe care incercam sa o minimizam/maximizam.

Vom folosi aceasta metoda pentru a incerca sa aproximam matricea A^{-1} , sau W , aceasta fiind matricea care ne va ajuta sa separam componentele, asa cum am mentionat pana acum. Functia de update pentru o abordare de tip gradient descent pentru optimizarea maximum likelihood este:

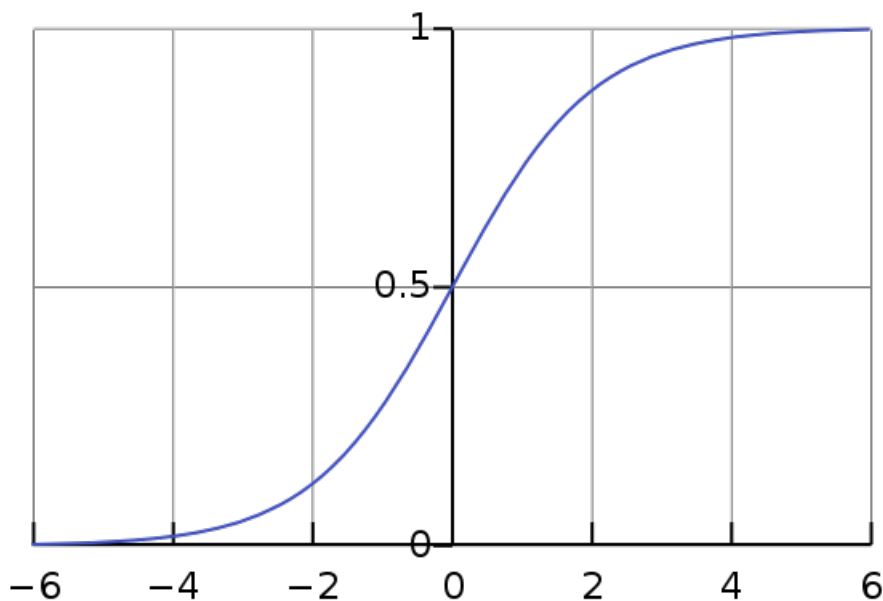
$$W^+ = W + \mu[I + g(y)y^T]W \quad (46)$$

unde μ este rata de invatare a algoritmului, iar g este functia de densitate a componentelor independente.

Un factor important in succesul acestui algoritm este alegerea functiei de densitate. Pana acum am presupus ca se cunoaste functia de densitate a componentelor independente, dar in practica, acest lucru nu poate fi garantat, iar din aceasta cauza, va fi nevoie de o aproximare. S-a observat ca o functie potrivita acestei probleme este functia sigmoida, deoarece avand o distributie continua, calculand derivata acestei functii ajungem la o functie de densitate, care poate fi folosita pentru aproximarea densitatii componentelor independente. Aici vom defini g ca fiind functia sigmoida, in practica s-a demonstrat ca este suficient de buna:

$$g_i(x) = \frac{1}{1 + e^{-x}} \quad (47)$$

Functia sigmoida are urmatoarea distributie:



4.9 Metodologia si algoritmul

4.9.1 Jupyter Notebook

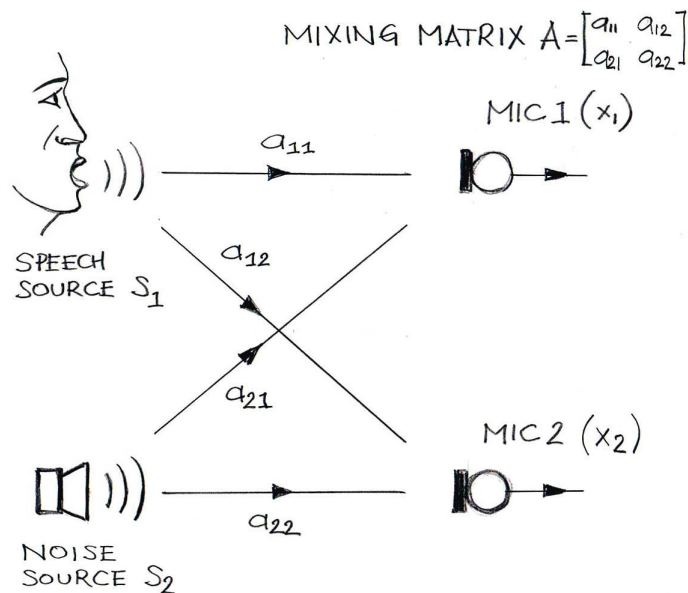
Implementarea algoritmului s-a facut in limbajul de programare Python, folosind platforma Jupyter Notebook. Am ales aceasta platforma deoarece face ca pasii algoritmului sa poata fi mai bine urmariti, iar afisarea rezultatelor sa se faca intr-o maniera usoara si concisa. Aplicatia functioneaza pe un server local, iar codul este scris intr-un *notebook (caiet)*, prezent intr-o pagina web, codul fiind impartit in *cells (celule)*, in fiecare celula putem scrie o bucata de cod, care mai apoi poate fi executata. Acest stil de lucru este de folos deoarece permite executarea fiecarei celule individual, scurtand timpul dedicat intretinerii caietului. De asemenea, in aceste celule putem amplasa grafice, histograme, elemente de tip *mark-up* pentru prezentarea unor informatii suplimentare, si elemente de tip audio sau video.

4.9.2 Semnalele audio testate

4.9.3 Pre-procesarea datelor

Pentru toate semnalele audio inafara de cele primitive, a fost nevoie sa centram datele, la fel ca in cazul PCA, pentru standardizare.

Pornind de la semnalele audio initiale, vom simula amestecarea lor, folosind operatorul de amestecare, A . In realitate, acest A ar fi dat de distantele dintre sursele de semnal si dispozitivele care le inregistreaza. In cazul nostru, A va fi o matrice patratica, fie de 2×2 , fie de 3×3 , in functie de caz. Acest proces arata ca mai inainte:



Pentru fiecare caz testat, am ales initial matricea A aleator, cu valorile distribuite uniform intre 0.01 si 0.1. Deoarece acest lucru a facut verificarea rezultatelor si modificarea algoritmului destul de greu, am ales dintre matricile aleatoare A cele care au demonstrat cel mai bine eficacitatea algoritmului. O astfel de matrice este:

$$\begin{bmatrix} 0.56804456 & 0.92559664 \\ 0.07103606 & 0.08712930 \end{bmatrix}$$

Vom prezenta mai tarziu rezultatele fiecarui test, iar acolo vor fi prezentate toate variabilele folosite pentru a se ajunge la rezultatele noastre.

4.9.4 Aplicarea algoritmului Gradient Descent

Avand sursele de semnal amestecate, trebuie sa aplicam algoritmul gradient descent prezentat mai devreme pentru a incerca sa gasim componentele independente initiale. Pentru acest lucru, vom porni de la o matrice A^{-1} , sau W , cu valori aleatoare uniform distribuite intre 0.01 si 0.1. Acest interval s-a observat ca fiind unul bun pentru cazurile noastre, dar el trebuie adaptat in functie de situatie.

Avand functia de actualizare a lui W , definita mai sus, putem incepe executarea algoritmului. In faza initiala, acesta se executa pentru un numar finit de pasi, ales in functie de situatie. Dupa terminarea executarii acestui numar de iteratii, W actualizat ar trebui sa fie cel care poate sa separe componentele. Acest lucru nu este intotdeauna adevarat, ceea ce inseamna ca ne trebuie o alta metoda de a opri cautarea solutiei optime.

4.9.5 Conditia de oprire

Algoritmii care aplica metoda gradient descent, de obicei sunt opriti atunci cand functia de cost ajunge sau este foarte aproape de minim. In general, se considera ca functia a ajuns la costul minim atunci cand diferentele dintre valorile consecutive obtinute de algoritm sunt mai mici decat un anumit *threshold* (*prag*). Alegere acestui prag este foarte importanta, deoarece un prag prea mic va duce la castiguri foarte mici de precizie dupa multe operatii, iar un prag mare poate va opri algoritmul inainte de a ajunge la valoarea minima. Cu toate ca aceasta idee pare destul de simpla, balansul dintre precizia rezultatului obtinut si viteza algoritmului este ceva ce trebuie masurat si ajustat de la caz la caz, mai ales in cazul functiilor complexe, precum a noastra.

4.9.6 Rezultate

5 Concluzii //TODO

6 Bibliografie //TODO

References

- [1] Aapo Hyvärinen and Erkki Oja. Independent component analysis: Algorithms and applications. 2000.

[1]