

Data Mining

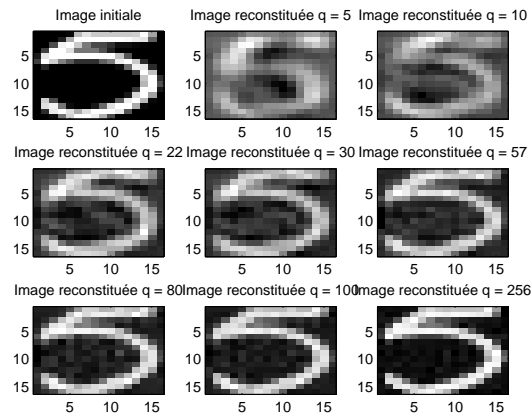
TP2 - Compression par ACP

Rémi MUSSARD - Thomas ROBERT

1 Tests de l'ACP

On utilise nos fonctions d'ACP sur des données contenant des images représentant des chiffres manuscrits. Après l'ACP, on reconstruit les données initiales.

```
1 load data_iris_usps_as1/uspsasi.mat
2
3 caract = 5;
4 lines = find(y == caract);
5
6 % récupération du caractère <caract>
7 X = x(lines, :);
8
9 % préparation de l'ACP
10 [D, U, moy]=mypca(X);
11
12 i = 2;
13
14 figure();
15 subplot(3,3,1);
16 imagesc(reshape(X(1,:),16,16)')
17 colormap gray;
18 title('Image initiale');
19
20 for q = [5 10 22 30 57 80 100 256]
21
22     % récupération des composantes principales
23     P = U(:, 1:q);
24
25     % projection sur les composantes principales
26     Ct = projpca(X, moy, P);
27
28     % reconstruction des données initiales
29     Xhat = reconstuctpca(Ct, P);
30
31     % representation de la reconstruction
32     subplot(3,3,i);
33     imagesc(reshape(Xhat(1,:),16,16)')
34     title(['Image reconstituée q = ' int2str(q)]);
35     colormap gray;
36
37     i = i + 1;
38 end
```



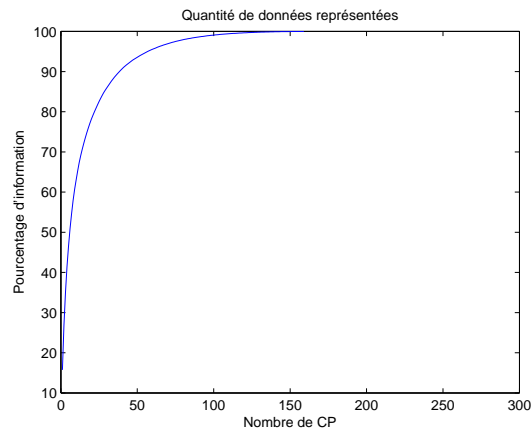
On remarque que comme prévu, plus on ajoute de composantes sur lesquelles on projette, plus l'image reconstituée est proche de l'image initiale. La différence entre l'image initiale et l'image avec toutes les composantes est sans doute due aux erreurs numériques lors des projections.

```

1 % nombre de CP approprié
2 figure();
3 qteRepresentee = cumsum(D)/sum(D)*100;
4 plot(qteRepresentee);
5 title('Quantité de données représentées');
6 xlabel('Nombre de CP');
7 ylabel('Pourcentage d''information');
8 fprintf('Nombre de CP pour avoir 95%% : %i\n', find(qteRepresentee > 95, 1));

```

Nombre de CP pour avoir 95% : 57



Avec 57, on aura 95% de l'information reconstruite. Etant donné l'allure de la courbe, cette valeur nous semble un bon compromis.

Visages propres

On charge les données et on calcule l'ACP sur le premier jeu.

```

1 % chargements des données
2
3 load YaleFace/Subset1YaleFaces.mat
4
5 X1 = X;
6 Y1 = Y;
7

```

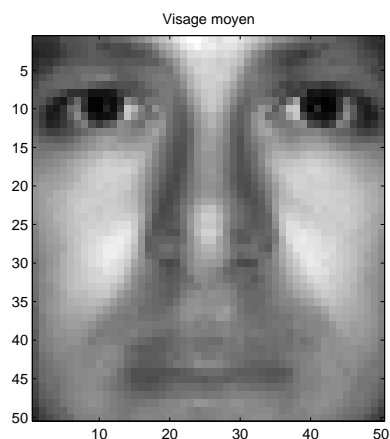
```
8 load YaleFace/Subset2YaleFaces.mat
9
10 X2 = X;
11 Y2 = Y;
12
13 load YaleFace/Subset3YaleFaces.mat
14
15 X3 = X;
16 Y3 = Y;
17
18 % calcul de l'ACP sur le premier ensemble
19 [D, U, moy]=mypca(X1);
```

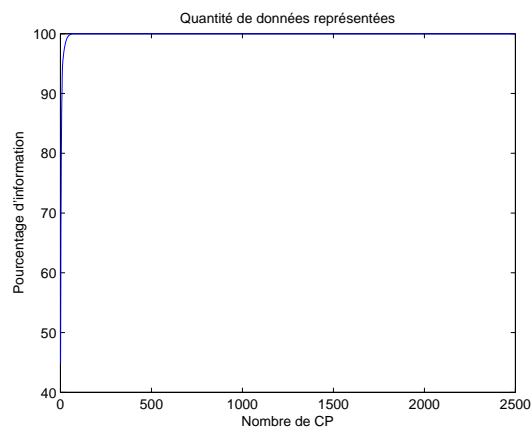
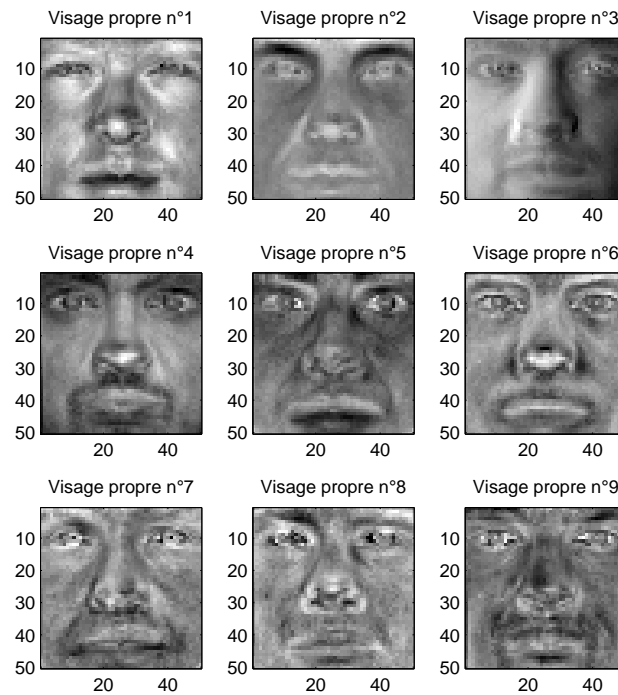
1.1 Visualisation des résultats de l'ACP

On affiche le visage moyen, les 9 premiers visages propres, et la quantité de données représentées.

```
1 % visage moyen
2 figure();
3 imagesc(reshape(moy,50,50));
4 colormap gray;
5 title('Visage moyen');
6
7 % visages propres
8 figure();
9 for i=1:9
10     subplot(3,3,i);
11     imagesc(reshape(U(:,i),50,50));
12     colormap gray;
13     title(['Visage propre n°' int2str(i)]);
14 end
15
16 % quantité d'info
17 figure();
18 qteRepresentee = cumsum(D)/sum(D)*100;
19 plot(qteRepresentee);
20 title('Quantité de données représentées');
21 xlabel('Nombre de CP');
22 ylabel('Pourcentage d''information');
23 fprintf('Nombre de CP pour avoir 95%% : %i\n', find(qteRepresentee > 95, 1));
```

Nombre de CP pour avoir 95% : 13





1.2 Reconstruction des données

Affichage des résultats de la compression par ACP.

```

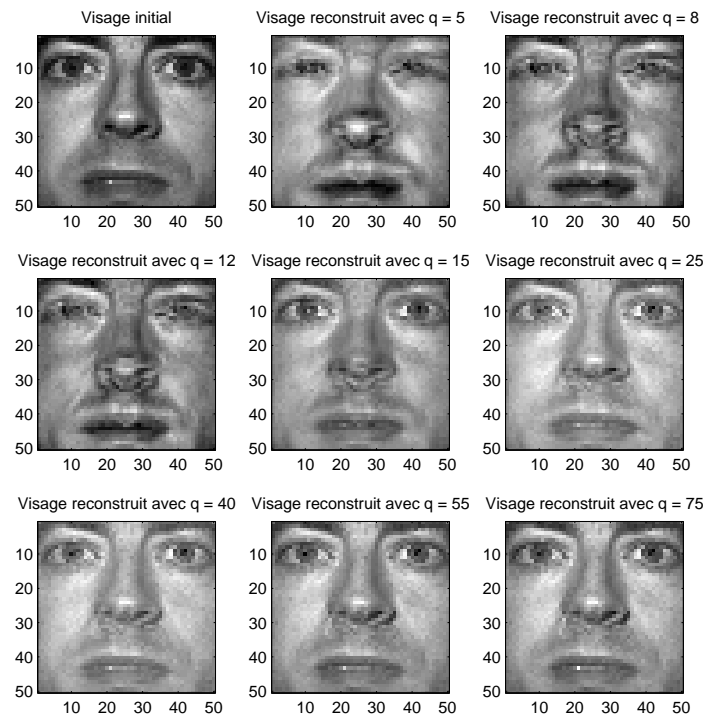
1 % visage a afficher
2 visage = 40;
3
4 i = 2;
5 figure();
6 subplot(3,3,1);
7 imagesc(reshape(X1(visage,:),50,50));
8 title('Visage initial');
9 colormap gray;
10
11 for q=[5 8 12 15 25 40 55 75]
12     % récupération des composantes principales
13     P = U(:, 1:q);
14
15     % projection sur les composantes principales
16     Ct = projpca(X1, moy, P);
17

```

```

18 % reconstruction des données initiales
19 Xhat = reconstuctpca(Ct, P);
20
21 % affichage
22 subplot(3,3,i);
23 imagesc(reshape(Xhat(visage,:),50,50));
24 title(['Visage reconstruit avec q = ' int2str(q)]);
25 colormap gray;
26
27 i = i + 1;
28 end

```



1.3 Détermination des meilleurs k et q

On détermine les meilleurs k et q en utilisant le premier jeu de données en apprentissage et le deuxième en validation.

```

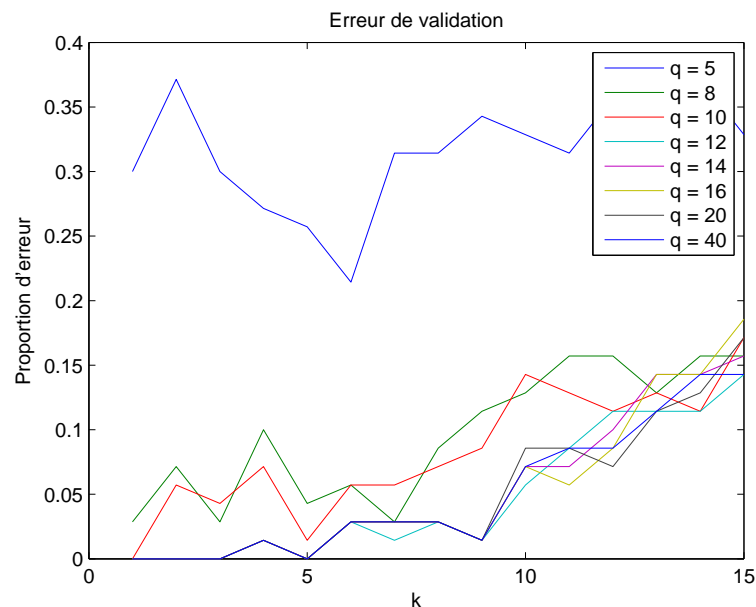
1 % pour chaque bloc
2 qVals = [5 8 10 12 14 16 20 40];
3 qValsLabel = {'q = 5', 'q = 8', 'q = 10', 'q = 12', 'q = 14', 'q = 16', 'q = 20', 'q = 40'};
4 kmax = 15;
5
6 errVal = zeros(kmax, length(qVals));
7 for j = 1:length(qVals)
8
9     % choix q
10    q = qVals(j);
11
12    % projection
13    P = U(:, 1:q);
14    Ct1 = projpca(X1, moy, P);
15    Ct2 = projpca(X2, moy, P);
16
17    % pour chaque k

```

```

18   for k = 1:kmax
19
20       % prédiction
21       [Y2pred, Dist] = knn(Ct2, Ct1, Y1, k);
22
23       % calcul proportion d'erreur pour le k choisi
24       errVal(k, j) = mean(Y2pred ~= Y2);
25   end
26 end
27
28 % affichage erreur
29 figure();
30 plot(errVal);
31 title('Erreur de validation');
32 xlabel('k');
33 ylabel('Proportion d''erreur');
34 legend(qValsLabel);

```



On voit qu'en choisissant $k = 5$, à partir d'une valeur suffisamment grande de q , l'erreur devient nulle. Ceci étant vrai à partir de $q = 12$, il n'est pas nécessaire de conserver plus de composantes qui alourdirait la base de données.

Les valeurs optimales sont donc $k = 5$ et $q = 12$.

1.4 Test de performance

On teste la performance de notre k -nn en utilisant le troisième jeu de données en test.

```

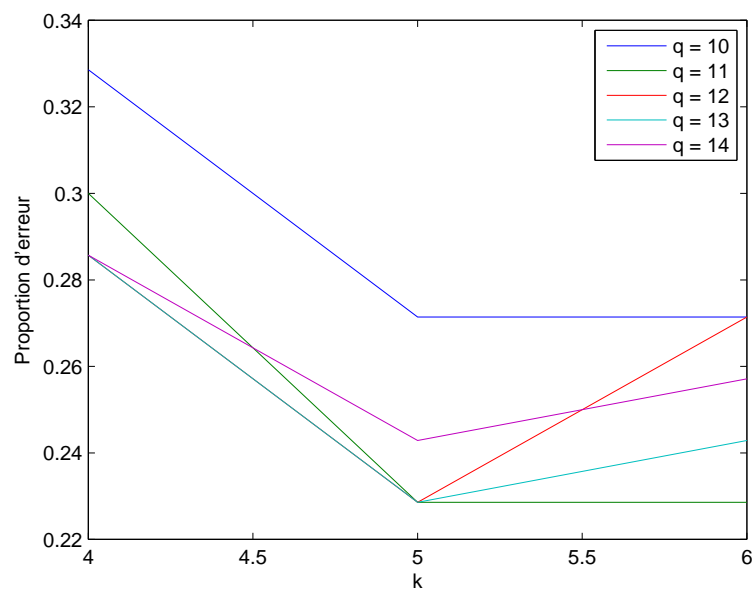
1 qTests = [10 11 12 13 14];
2 qTestsLabel = {'q = 10', 'q = 11', 'q = 12', 'q = 13', 'q = 14'};
3 kTests = 4:6;
4
5 errTest = zeros(length(kTests), length(qTests));
6 for j = 1:length(qTests)
7
8     % choix q
9     q = qTests(j);
10
11     % projection
12     P = U(:, 1:q);
13     Ct1 = projpca(X1, moy, P);
14     Ct3 = projpca(X3, moy, P);

```

```

15
16 % pour chaque k
17 for i = 1:length(kTests)
18
19     % choix k
20     k = kTests(i);
21
22     % prédiction
23     [Y3pred, Dist] = knn(Ct3, Ct1, Y1, k);
24
25     % calcul proportion d'erreur pour le k choisi
26     errTest(i, j) = mean(Y3pred ~= Y3);
27 end
28 end
29
30 % affichage erreur
31 figure();
32 plot(kTests, errTest);
33 xlabel('k');
34 ylabel('Proportion d''erreur');
35 legend(qTestsLabel);

```



On voit que les valeurs $k = 5$ et $q = 12$ sont effectivement les valeurs optimales pour les résultats sur l'ensemble de test. La proportion d'erreur reste de 22% mais c'est le meilleur résultat qu'on ait pu obtenir sur l'ensemble de test. Afin d'augmenter ce résultat, il faudrait agrandir la base d'apprentissage.

2 Conclusion

Ce TP aura permis de voir que l'ACP permet de compresser les données de façon significative. Par exemple, il est possible de conserver quasiment l'intégralité des données d'une photo sur 13 composantes au lieu de 2500.

Cette compression permet par ailleurs de pouvoir utiliser des méthodes comme le k -nn dans une base de taille réduite par rapport aux données initiales et donc d'augmenter grandement la rapidité des calculs.