

# Podjęcia Systemowe w Badaniach Biomedycznych

## zajęcia laboratoryjne 2

### Redukcja wymiarowości

## 1 Redukcja wymiarowości

Redukcja wymiarowości jako ważny element eksploracji danych jest to proces transformacji danych wielowymiarowych do przestrzeni o mniejszym wymiarze. Podejście to jest stosowane ze względów praktycznych, celem ułatwienia wizualizacji danych. Najczęściej dane są redukowane są do 2 lub 3 wymiarów. Dodatkowo, po zmniejszeniu wymiaru, wiele zależności wydaje się być bardziej czytelnych.

Celem metod redukcji danych jest zidentyfikowanie ukrytego wzorca danych, zmniejszenie wymiarowości poprzez usunięcie szumów i redundantnych danych oraz zidentyfikowanie skorelowanych zmiennych.

W R dostępnych jest wiele różnych funkcji i pakietów dla przeprowadzenia redukcji wymiarowości różnymi metodami. Wyniki mogą się różnić pomiędzy tymi pakietami. W tym protokole skupimy się tylko na wybranych z nich.

## 2 Ocena danych wejściowych

Często otrzymane dane wejściowe, bez wcześniejszego przetwarzania nie nadają się do bezpośredniego użycia i przeprowadzenia analiz. W danych mogą pojawiać się braki i wartości odstające, które mogą wpływać na wyniki analiz. Opis tych właściwości i sposoby radzenia sobie z nimi zostały omówione na poprzednich zajęciach i wykładzie.

Należy jednak wspomnieć także o normalizacji i skalowaniu danych. Konieczność takiej obróbki wynika z różnego charakteru danych. Jeśli zmienne charakteryzują się różną skalą to ważne jest zapewnienie porównywalności tych zmiennych. R w niektórych funkcjach domyślnie wykonuje takie przetwarzanie danych. Wyróżniamy następujące przekształcenia danych:

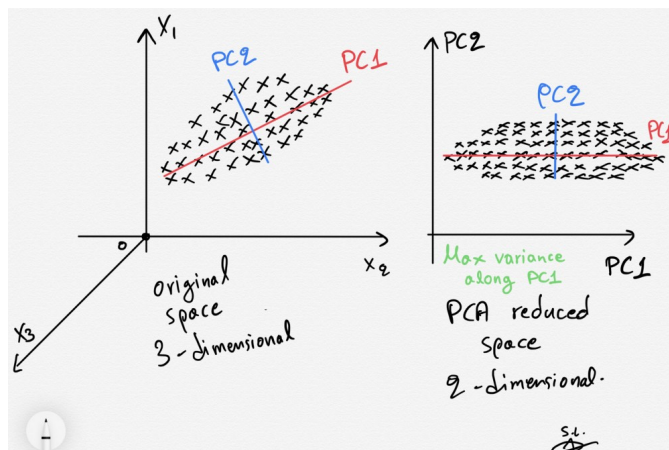
- normalizacja min-max: dane po normalizacji należą do przedziału od 0 do 1,
- standaryzacja: dane po standaryzacji będą miały średnią = 0 i odchylenie standardowe = 1,
- centrowanie: centrowanie danych wokół wartości 0.

## 3 Analiza składowych głównych - PCA

Analiza składowych głównych (*ang. Principal Components Analysis - PCA*) pozwala podsumować i wizualizować informacje w zbiorze danych, które zawierają obserwacje opisane przez wiele wzajemnie skorelowanych zmiennych ilościowych (każda zmienna może zostać uznana za inny wymiar). Zastanów się, jak wizualizować dane jeśli mamy więcej niż trzy zmienne? Analiza składowych głównych służy do wyodrębnienia ważnych informacji z wielowymiarowych danych i wyrażenia tych informacji jako zestawu kilku nowych zmiennych zwanych właśnie składowymi głównymi (*ang. Principal Components - PC*), takie nowe zmienne odpowiadają liniowej kombinacji oryginałów. Celem analizy PCA jest znalezienie transformacji układu współrzędnych, która lepiej opisze zmienność pomiędzy obserwacjami. Innymi słowy, chodzi o zidentyfikowanie kierunków (składowych głównych), wzdłuż których zmienność danych jest maksymalna. Przy redukcji wymiarowości danych zależy nam na ograniczeniu ilości wielowymiarowych danych do dwóch lub trzech składowych głównych, przy minimalnej utracie informacji.

### 3.1 Zrozumieć PCA

Zrozumienie szczegółów PCA wymaga znajomości algebry liniowej. Najprościej jednak można wyjaśnić podstawy na bazie prostej graficznej reprezentacji danych. Na poniższym wykresie 1 po lewej stronie, dane są przedstawione w przestrzeni 3D. Zmniejszenie wymiarów uzyskuje się poprzez określenie głównych kierunków, zwanych składowymi głównymi, w których dane się zmieniają. Na rysunku 1 po lewej stronie, oś PC1 jest pierwszym głównym kierunkiem, wzdłuż którego próbki wykazują największą zmienność. Oś PC2 jest drugim najważniejszym kierunkiem i jest prostopadła do osi PC1. Wymiarowość 3D można zredukować do wymiaru 2D, co przedstawia wykres 1 po prawej stronie.



Rysunek 1: Zrozumienie koncepcji PCA na przykładzie redukcji wymiaru 3D do 2D.

Źródło: <https://ai.plainenglish.io/how-to-implement-pca-with-python-and-scikit-learn-22f3de4e5983>

Jak to rozumieć od strony technicznej? Za pomocą pewnych transformacji liniowych, tworzony jest nowy układ współrzędnych, gdzie pierwsza oś ma największą wariancję, a druga oś ma drugą w kolejności największą wariancję. Składowe, których wariancja jest niska uznane są za nieważne i można je pominąć. Nowo utworzony zbiór zmiennych, ma być odpowiednio mały, ale tak dobrany aby odzwierciedlić zmienność występującą w oryginalnym zestawie danych. Wielkość wariancji zachowana przez każdą składową główną jest mierzona tak zwaną wartością własną (ang. *eigenvalue*).

### 3.2 Przykład *prcomp()* i *princomp()* w R

Najbardziej podstawowymi funkcjami do analizy PCA jest *prcomp()* i *princomp()* z pakietu *stats*, jednak każda z tych funkcji wykorzystuje inną metodę:

- *prcomp()* - wyznacza składowe główne przy użyciu dekompozycji SVD,
- *princomp()* - wyznacza składowe główne poprzez wektory własne macierzy kowariancji/korelacji.

W obu powyższych przypadkach należy podać dane wejściowe i określić pozostałe argumenty. W przypadku funkcji *prcomp* należy określić czy dane mają zostać przeskalowane przed przystąpieniem do analizy, służy do tego argument *scale* (*scale = TRUE* jeśli dane mają być skalowanie lub *scale = FALSE* jeśli nie). Podobnie jest w przypadku funkcji *princomp()*, ustawienie argumentu *cor = TRUE* określa, że dane mają zostać przeskalowane i centrowane przed analizą. Opisane wyżej metody teoretycznie dają identyczne wyniki (właściwości składowych głównych będą identyczne), jednak w niektórych sytuacjach wyniki mogą się różnić.

Dane wykorzystywane do poniższych danych pochodzą z pakietu *factoextra*:

```
1 help(package="factoextra") #zobacz co oferuje pakiet
2 install.packages("factoextra")
3 library("factoextra")
4
```

```
5 data(decathlon2) #dane z pakietu factoextra
6 dane <- decathlon2[1:10]
```

---

Przykład użycia funkcji *prcomp()*:

---

```
1 #analiza PCA
2 result.pcomp <- prcomp(dane, scale = TRUE, center = TRUE)
3 result.pcomp
4
5 #wykres obserwacji (individuals/x/scores)
6 result.pcomp$x
7 plot(result.pcomp$x, main = "scale = TRUE & center = TRUE")
8 text(result.pcomp$x, labels = row.names(dane), pos = 3, cex = 0.7)
9
10 #wykres zmiennych (variables/rotation/loadings)
11 rotation <- result.pcomp$rotation
12 rotation
13 plot(rotation, main = "scale = TRUE & center = TRUE")
14 text(rotation, labels = colnames(dane), pos = 2, cex = 0.7)
15 for(i in 1:dim(dane)[2]) {
16   arrows(0, 0, rotation[i, 1], rotation[i, 2], angle = 7, col = "red")
17 }
18
19 #wykres obserwacji i zmiennych
20 biplot(result.pcomp)
21
22 #wykres osuwiskowy – wykres wariacji
23 screeplot(result.pcomp, main = "scale = TRUE & center = TRUE", type = "l")
24
25 #wykres osuwiskowy można także otrzymać przy użyciu funkcji plot()
26 plot(result.pcomp, main = "scale = TRUE & center = TRUE")
27
28 #podsumowanie danych przedstawionych na wykresie screeplot() i plot()
29 summary(result.pcomp)
```

---

Przykład użycia funkcji *princomp()*:

---

```
1 #analiza PCA: cor = T macierz korelacji, cor = F macierz kowariancji
2 result.princomp <- princomp(dane, cor = TRUE)
3 result.princomp
4
5 #wykres obserwacji (individuals/x/scores)
6 result.princomp$scores
7 plot(result.princomp$scores)
8 text(result.princomp$scores, labels = row.names(dane), pos = 3, cex = 0.7)
9
10 #wykres zmiennych (variables/rotation/loadings)
11 loading <- result.princomp$loadings
12 loading
13 plot(loading)
14 text(loading, labels = colnames(dane), pos = 2, cex = 0.7)
15 for(i in 1:dim(dane)[2]) {
16   arrows(0, 0, loading[i, 1], loading[i, 2], angle = 7, col = "red")
17 }
18
```

```

19 #wykres obserwacji i zmiennych
20 biplot(result.princomp)
21
22 #wykres osuwiskowy
23 screeplot(result.princomp, type = "l")
24
25 #wykres osuwiskowy mozna takze otrzymac przy uzyciu funkcji plot()
26 plot(result.princomp)
27
28 #podsumowanie danych przedstawionych na wykresie screeplot() i plot()
29 summary(result.princomp)

```

---

W kontekście wyników uzyskanych z funkcji *prcomp* i *princomp* należy zwrócić uwagę na inne nazwy obiektów, które w rzeczywistości wskazują te same wyniki:

- zmienna *x* dla *prcomp()* i zmienna *scores* dla *princomp()*: współrzędne obserwacji składowych głównych,
- zmienna *rotation* dla *prcomp()* i zmienna *loadings* dla *princomp()*: wektory ładunków.

### 3.3 Pakiet *factoextra*

Pakiet *factoextra* został stworzony by umożliwiać wydobywanie danych i lepszą wizualizację dla różnych metod redukcji wymiaru, poniżej zapoznamy się z wybranymi funkcjami.

Wydobywanie danych:

---

```

1 #wyniki uzyskane z wcześniejszej analizy
2 result.pcomp
3
4 #otrzymanie wartosci wlasnych i wariancji składowych glownych
5 get_eigenvalue(result.pcomp)
6
7 #otrzymanie wyników dla obserwacji
8 get_pca_ind(result.pcomp)
9 #Principal Component Analysis Results for individuals
10 #=====
11 # Name      Description
12 #1 "$coord"  "Coordinates for the individuals"
13 #2 "$cos2"   "Cos2 for the individuals"
14 #3 "$contrib" "contributions of the individuals"
15 get_pca_var(result.pcomp)$coord
16 get_pca_var(result.pcomp)$cos2 #jakosc reprezentacji
17 get_pca_var(result.pcomp)$contrib #udzial w składowych glownych
18
19 #otrzymanie wyników dla zmiennych
20 get_pca_var(result.pcomp)

```

---

Jaka jest interpretacja powyższych wyników?

- **Wartości własne i wariancje składowych głównych:** w wyniku działania funkcji *get\_eigenvalue()*, otrzymujemy tabelę z trzema kolumnami (*eigenvalue*, *variance.percent*, *cumulative.variance.percent*). Pierwsza kolumna zawiera wartości własne, suma wszystkich tych wartości daje całkowitą wariancję równą 10. Druga kolumna zawiera procent wariancji, który oznacza że  $X_1\%$  zmienności jest wyjaśnione przez konkretną wartość własną. Trzecia kolumna zawiera skumulowany procent wariancji (uzyskany przez sumę kolejnych procentów wariancji ( $X_1\% + X_2\% + \dots$ )). Załóżmy, że mowa o sumie dwóch pierwszych procentów wariancji (dla 2 pierwszych wartości własnych), zatem suma  $X_1\% + X_2\%$  jest to procent zmienności wyjaśniony przez dwie pierwsze wartości własne (wartości własne określają poszczególne składowe główne).

Sprawdzamy wartości własne, aby określić liczbę składowych głównych, które mają zostać zachowane po analizie PCA. Wybierane są te, które mają największe wartości, co ma na celu minimalizację straty informacji podczas redukcji wymiarów. Wartości własne są wyższe dla pierwszych składowych głównych i niższe dla kolejnych. Wartości własne są miarą zmienności pierwotnych danych przedstawionych we współrzędnych składowych głównych.

- **Wyniki dla obserwacji i zmiennych:** funkcja `get_pca_ind()` wydobywa dane dla obserwacji, podczas gdy funkcja `get_pca_var()` wydobywa informacje dla zmiennych. Poniżej znajduje się lista obiektów, przez które możemy dostać się do szczegółowych danych:
  - `coord` - wydobywa współrzędne zmiennych wykorzystywane w celu utworzenia wykresu punktowego,
  - `cor` - wydobywa dane o korelacjach między zmiennymi a wymiarami,
  - `cos2` - wydobywa dane o jakości reprezentacji zmiennych,
  - `contrib` - wydobywa dane o udziale procentowym zmiennych w składowych głównych.

Uzyskane dane wykorzystuje się na dalszym etapie do wizualizacji.

Wizualizacja danych:

---

```
1 #Wizualizacja wartosci wlasnych na wykresie osuwiskowym
2 fviz_eig(result.pcomp)
3 fviz_eig(result.pcomp, addlabels = TRUE)
4
5 #wykres obserwacji (individuals/x/scores)
6 fviz_pca_ind(result.pcomp)
7 #dodatkowe mozliwosci reprezentacji obserwacji
8 fviz_pca_ind(result.pcomp,
9   col.ind = "cos2", #kolor wedlug jakosci reprezentacji
10  repel = TRUE #nie nakladaja sie nazwy
11 )
12
13 #wykres jakosci reprezentacji
14 fviz_cos2(result.pcomp, choice = "ind", axes = 1:2) #2 wymiary
15 #fviz_cos2(result.pcomp, choice = "ind") # 1 wymiar
16
17 #wykres udzialu w skladowych glownych
18 fviz_contrib(result.pcomp, choice = "ind", axes = 1:2) #2 wymiary
19
20 #wykres zmiennych (variables/ratation/loadings)
21 fviz_pca_var(result.pcomp)
22 #dodatkowe mozliwosci reprezentacji obserwacji
23 fviz_pca_var(result.pcomp,
24   col.var = "contrib", #kolor wedlug udzialu w skladowych glownych
25   repel = TRUE
26 )
27
28 #wykres jakosci reprezentacji
29 fviz_cos2(result.pcomp, choice = "var", axes = 1:2) #2 wymiary
30
31 #wykres udzialu w skladowych glownych
32 fviz_contrib(result.pcomp, choice = "var", axes = 1:2) #2 wymiary
33
34 #biplot
35 fviz_pca_biplot(result.pcomp, repel = TRUE)
```

---

Jaka jest interpretacja powyższych wyników?

- **Wykres osuwiskowy:** wykres osuwiskowy utworzony za pomocą funkcji `fviz_eig()` przedstawia procent wariancji dla każdego wymiaru. Można użyć dodatkowego argumentu `addlabels = TRUE`, aby dodać etykiety z dokładnymi wartościami. Wykres ten jest graficzną wizualizacją zmienności, można określić ile procent zmienności pierwotnych danych wyjaśnia poszczególna składowa główna.
- **Wykres dla zmiennych:** podstawowy wykres dla zmiennych można utworzyć funkcją `fviz_pca_var()`. Wykres dla zmiennych pokazuje wpływ zmiennych na poszczególne składowe główne. Zmienne są reprezentowane przez wektory (zwane wektorami ładunków), a ich kierunek i długość określają w jakim stopniu wpływają na poszczególne składowe. Gdy zmienne leżą blisko siebie są silnie dodatnio skorelowane, gdy leżą po przeciwnych stronach są silnie ujemnie skorelowane, a gdy leżą względem siebie prostopadłe to nie wykazują żadnej korelacji. Ponadto, w przypadku zmiennych na wykresie może pojawić się tzw. koło korelacji: im większa wartość to zmienne ułożone są bliżej obwodu koła korelacji, a im niższa tym zmienne ułożone są bliżej środka tego koła i są mniej ważne dla składowych głównych - odległość zmiennych jest miarą jakości.

Dla wykresu zmiennych można użyć dodatkowych argumentów jak `col.var`, żeby otrzymać dodatkową informację kolorystyczną, np. o jakości reprezentacji (`col.var = "cos2"`) czy o udziale w składowych głównych (`col.var = "contrib"`). Wspomniane dane, można narysować także na oddzielnych wykresach:

- wykres jakości reprezentacji `fviz_cos2()` z argumentem `choice = "var"`: wysoka wartość świadczy o dobrej reprezentacji zmiennej w składowej głównej.
- wykres udziału w składowych głównych `fviz_contrib()` z argumentem `choice = "var"`: im większa wartość, tym większy udział zmiennej w składowej głównej. Czerwona przerywana linia na tym wykresie wskazuje oczekiwany średni wkład, zmienne przekraczające ten próg mają największy udział w badanych składowych głównych.
- **Wykres dla obserwacji:** podstawowy wykres dla obserwacji można utworzyć funkcją `fviz_pca_ind()`. Wykres ten pokazuje, które grupy obserwacji są do siebie podobne (są zgrupowane razem), a które się od siebie różnią (są od siebie oddalone). Podobnie jak dla wykresu zmiennych, również dla wykresu obserwacji można użyć dodatkowych argumentów jak `col.ind`, żeby otrzymać dodatkową informację kolorystyczną, np o jakości reprezentacji (`col.ind = "cos2"`) czy o udziale w składowych głównych (`col.ind = "contrib"`). Wspomniane dane, można narysować także na oddzielnych wykresach (podobnie jak dla zmiennych, należy tylko ustawić odpowiednio argumenty):
  - wykres jakości reprezentacji `fviz_cos2()` z argumentem `choice = "ind"`,
  - wykres udziału w składowych głównych `fviz_contrib()` z argumentem `choice = "ind"`.
- **Wykres obserwacji i zmiennych:** funkcja `biplot()` należy zdawać sobie sprawę, że obserwacje i zmienne nie leżą w tej samej przestrzeni, dlatego powinno się skupiać na kierunku, a nie na pozycjach bezwzględnych na wykresie. Im dana obserwacja leży bliżej od danej zmiennej, tym większy jest wpływ tej zmiennej na daną obserwację.

## Zadania

Sprawozdanie i kod R prześlij na maila (Kaja.Gutowska@cs.put.poznan.pl). Proszę nadać tytuł maila zgodnie z opisem: [P\\$wBB\\_Imię\\_Nazwisko\\_lab2](#).

### Zad. 1

Przetestuj funkcję `prcomp()` i `princomp()`. Zaobserwuj jakie są wyniki przy zastosowaniu skalowania i centrowania oraz jak wyglądają dane gdy nie zastosujemy transformacji danych. Pamiętaj, że obie te metody używają innych argumentów. Jakie wartości, przy jakich argumentach pozwalają uzyskać takie same wyniki i wykresy? Jakie są konsekwencje, gdy nie przeprowadzimy skalowania danych?

### Zad. 2

Wybierz i wykorzystaj istniejące dane np. ze strony Głównego Urzędu Statystycznego (<https://stat.gov.pl>) lub dowolne inne dane, aby przeprowadzić redukcję wymiarowości metodą PCA. Opisz wybrane dane, poszczególne etapy analizy i uzyskane wyniki.