

Analizne metode za karakterizacijo materialov in bioloških sistemov

izr.prof.dr. Helena Prosen
izr.prof.dr. Matevž Pompe

Metoda glavnih osi (Principal Component Analysis- PCA)

Uvod

Metoda glavnih osi (PCA) je linearna transformacija (rotacija) m -dimenzionalnega merskega prostora. Prva koordinata zasukanega koordinatnega sistema (glavna os ali 1. PC) je tista smer, v kateri se celotna varianca sistema vseh obravnavanih podatkov najbolj zmanjša. Naslednja glavna os je pravokotna na prejšnjo in je zopet določena s tem, da opiše kar največ preostale variance, itd. Ker je tako velika večina celotne variance podatkov zbrana okrog nekaj prvih novih osi, lahko ostale osi zanemarimo, posebej, če upoštevamo samo tiste, ki so skupno odgovorne za več kot rpnj 90 % variance. Ker so v veliki večini kemijskih problemov spreminjivke med seboj odvisne, lahko že prvi dve glavni (novi) osi nosita več kot 75 % variance. Glavne osi so linearne kombinacije starih spremenljivk:

$$PC_1 = l_{11}x_1 + l_{12}x_2 + \dots + l_{1n}x_n = \sum_{i=1}^n l_{1i}x_i \rightarrow \text{prva glavna os}$$

$$\text{splošno: } PC_j = \sum_{i=1}^n l_{ji}x_i \rightarrow j\text{-ta glavna os}$$

$$\sum_{i=1}^n l_{ji}^2 = 1 \rightarrow \text{za vsako os } PC_j \quad \sum_{j=1}^m l_{ji} = 0 \rightarrow \text{za vsak poljubni par } PC_i \text{ in } PC_j$$

Koeficienti l_{ji} (*loadings*) povedo, koliko je vsaka originalna (stara) spremenljivka x_i udeležena v posamezni novi osi PC_j .

Ker je večina informacije (variance) skoncentrirana okoli prvih dveh novih osi, lahko PCA služi kot projekcija objektov iz m -dimenzionalnega prostora v 2d prostor prvih dveh glavnih osi (projekcijo originalnih objektov lahko naredimo v 2d prostor katerikoli drugih dveh novih osi).

Izhodišče za izračun glavnih osi oziroma koordinatnega sistema v katerem so osi po vrsti usmerjene tako, da pojasnijo kar največ variance, je *matrica korelacij* med vsemi spremenljivkami. Če med spremenljivkami ni korelacij, torej v primeru, ko imamo opraviti z med seboj povsem neodvisnimi spremenljivkami, glavnih osi ne moremo izračunati!

2

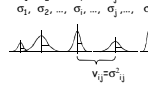
Podatki za izračun glavnih osi

Čeprav ima danes večina statističnih programov, zlasti tistih, ki so namenjeni kemikom, metodo PCA že vgrajeno, si vseeno oglejmo matematični postopek izračuna glavnih osi. Določitev glavnih osi pričnemo z $n \times m$ veliko matriko, ki je izračunana na podlagi n izmerjenih podatkov X . Posamezni objekt k je predstavljen v m -dimenzionalnem prostoru kot vektor $X_k = [x_{k1}, x_{k2}, \dots, x_{km}]$. Za izračun PCA lahko matriko podatkov X predhodno preoblikujemo na tri načine. Preoblikovano matriko bomo pisali kot matriko A_i :

$$\left. \begin{array}{l} X_1 \rightarrow x_{11}, x_{12}, \dots, x_{1n}, \dots, x_{1m} \\ X_2 \rightarrow x_{21}, x_{22}, \dots, x_{2n}, \dots, x_{2m} \\ \vdots \\ X_k \rightarrow x_{k1}, x_{k2}, \dots, x_{kn}, \dots, x_{km} \\ \vdots \\ X_n \rightarrow x_{n1}, x_{n2}, \dots, x_{nn}, \dots, x_{nm} \\ \bar{X} \rightarrow \bar{x}_1, \bar{x}_2, \dots, \bar{x}_1, \dots, \bar{x}_m \\ \sigma_1, \sigma_2, \dots, \sigma_1, \dots, \sigma_m \end{array} \right\} \begin{array}{l} A_1 = \|x_{ij}\| \rightarrow \text{originalni podatki} \\ A_2 = \|x_{ij} - \bar{x}_j\| \rightarrow \text{podatki premaknjeni za povprečje} \\ A_3 = \left\| \frac{x_{ij} - \bar{x}_j}{\sigma_j} \right\| \rightarrow \text{normalizirani podatki } (\bar{x}_j = 0, \sigma_j = 1) \end{array}$$

$$\text{varianca} \rightarrow v_j = \frac{\sum_{k=1}^n (x_{kj} - \bar{x}_j)^2}{m-1}$$

$$\text{ko varianca} \rightarrow v_j = \frac{\sum_{k=1}^n (x_{kj} - \bar{x}_j)(x_{kj} - \bar{x}_j)}{m-1}$$

$$\text{korelacija} \rightarrow r_{ij} = \frac{v_{ij}}{\sqrt{v_i v_j}} = \frac{\sum_{k=1}^n (x_{ki} - \bar{x}_i)(x_{kj} - \bar{x}_j)}{\sqrt{\sum_{k=1}^n (x_{ki} - \bar{x}_i)^2 \sum_{k=1}^n (x_{kj} - \bar{x}_j)^2}}$$


Z množenjem vsake od omenjenih treh matrik A_1 , A_2 , ali A_3 same s seboj v transponirani obliki (indeks 1!), dobimo tri različne izhodne matrike. Iz produkta $A_1 A_1^T$ lahko z diagonalizacijo izračunamo glavne osi PC, ($i = 1 \dots m$).

$X = (A_1)^T A_1$ je matrika mešanih produktov faktorjev
 $V = (A_2)^T A_2$ je matrika varianc in kovarianc
 $R = (A_3)^T A_3$ je korelacijska matrika

3

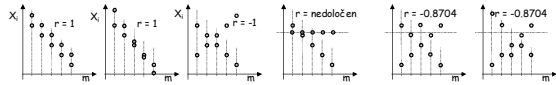
Izrazi za kovarianco in korelacijo med dvema nizoma vrednosti X_i in X_j

$$\text{kovarianca } v_{ij} = \frac{\sum_{k=1}^m (x_{ik} - \bar{x}_i)(x_{jk} - \bar{x}_j)}{m-1} = \frac{\sum_{k=1}^m x_{ik}x_{jk} - \frac{1}{m}(\sum_{k=1}^m x_{ik})(\sum_{k=1}^m x_{jk})}{m-1}$$

$$\text{korelacija } r_{ij} = \frac{\sum_{k=1}^m (x_{ik} - \bar{x}_i)(x_{jk} - \bar{x}_j)}{\sqrt{\sum_{k=1}^m (x_{ik} - \bar{x}_i)^2} \sqrt{\sum_{k=1}^m (x_{jk} - \bar{x}_j)^2}} = \frac{\sum_{k=1}^m x_{ik}x_{jk} - \frac{1}{m}(\sum_{k=1}^m x_{ik})(\sum_{k=1}^m x_{jk})}{\sqrt{\sum_{k=1}^m x_{ik}^2 - \frac{1}{m}(\sum_{k=1}^m x_{ik})^2} \sqrt{\sum_{k=1}^m x_{jk}^2 - \frac{1}{m}(\sum_{k=1}^m x_{jk})^2}}$$

Kovarianca v_{ij} in korelacija r_{ij} med dvema enako dolgima nizoma vrednosti: $X_i = (x_{i1}, x_{i2}, \dots, x_{im}, \dots, x_{im})$ in $X_j = (x_{j1}, x_{j2}, \dots, x_{jm}, \dots, x_{jm})$ se v statistiki računata dokaj pogosto. Obe vrednosti kažeta na medsebojno odvisnost obeh nizov števil. Veliko pogosteje kot kovarianca v_{ij} se uporablja korelacijski koeficient r_{ij} , ker je normaliziran in kot tak primerljiv med različnimi pari enako dolgih nizov.

Korelacijski koeficient r lahko zavzame vse realne vrednosti med -1 in 1. Če je $r = 0$, potem sta oba vektorja X_i in X_j med seboj neodvisna, če pa je $|r| = 1$, potem sta oba vektorja povezana z linearno funkcijo (enega lahko izračunamo iz drugega preko linearne enačbe): $X_j = aX_i + b$.



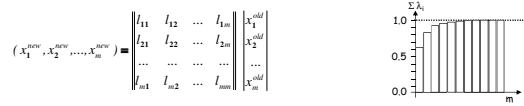
4

Izračun glavnih osi

Izračun glavnih osi je rotacija koordinatnega sistema. Rezultat PCA je transformacijska (rotacijska) matrika L , ki omogoča predstavitev objektov v novem koordinatnem sistemu. Matrika $||L||$ definira novi sistem tako, da je velikost deleža od celotne informacije (odstotek od celotne variance sistema vseh objektov) razporejen zaporedno po osih: prva os prevzame največji, zadnja os pa najmanjši delež celotne variance (informacije).

$$(X^{new})^T = ||L|| X^{old} \quad X_{ik}^{new} = \sum_{j=1}^m l_{jk} X_{ij}^{old}$$

Matriko $||L||$ dobimo z diagonalizacijo ene od treh matrik, omenjenih na prejšnji strani: $||X||$ (matrike mešanih produktov), $||V||$ (matrike varianc) ali $||R||$ (matrike korelacij). Z diagonalizacijo katerekoli od omenjenih treh $m \times m$ dimensionalnih matrik, dobimo najprej m lastnih vrednosti $\lambda_1, \lambda_2, \dots, \lambda_m$. Vsaki od lastnih vrednosti λ_i pripada lastni vektor $l_i = (l_{i1}, l_{i2}, \dots, l_{im})$. Koeficiente l_{ij} imenujemo komponente lastnih vektorjev (loadings). Ko lastne vrednosti λ_i uredimo po velikosti in v istem vrstnem redu razporedimo po vrsticah tudi lastne vektorje, dobimo transformacijsko matriko $||L||$ za izračun glavnih osi.



Objekti X^{new} , ki so predstavljeni v koordinatnem sistemu prvih dveh glavnih osi (PC_1 in PC_2), so dejansko predstavljeni samo s prvima dvema novima komponentama $X^{new} = (x_1^{new}, x_2^{new}, 0, 0, \dots, 0)$. Kolikor večji odstotek celotne informacije nosita prvi dve novi komponenti, toliko bolj verodostojna je preslika objektov X_i v 2-dimenzionalni prostor PC_1/PC_2 .

5

Primer: PCA analiza oljčnih olj

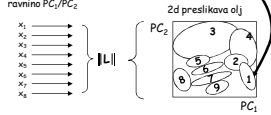
V danem primeru bomo s PCA obdelali 572 oljčnih olj, od katerih je vsako analizirano glede na vsebnost osmih maščobnih kislin. V spodnji tabeli so podani osnovni podatki: povprečje posameznih spremenljivk (koncentracije maščobnih kislin v 100 ppm), standardni odklik za vsako maščobno kislino in prvi dve vrstici izhodnih treh matrik $||X||$, $||V||$ in $||R||$. Vse tri matrike imajo dimenzijo 572 x 8.

	Palmitic-Palmitoleic	Stearic	Oleic	Linoleic	Arachidic	Linolenic	Eicosenic		
S. Apulija	1075	75	226	7823	672	36	60	29	X1
Kalabrija	1355	139	230	7299	832	42	60	32	X2
S. Apulija	-56.8	-51.1	-2.9	511.3	-308.5	4.0	1.8	12.7	V1-X1,XP
Kalabrija	83.3	32.9	1.1	-12.7	-148.5	30.0	1.8	15.7	V2-X2,XP
S. Apulija	-0.93	-0.97	-0.08	1.26	-1.27	0.32	0.08	0.90	R1-V1/s
Kalabrija	0.49	0.25	0.03	-0.03	-0.61	0.79	0.08	1.12	R2-V2/s
Povprečje	1231.8	126.1	228.9	7311.7	980.5	32.0	58.2	16.3	XP
std. odklik	568.6	52.5	36.7	405.8	242.8	12.7	21.8	14.1	s

Analyze 572 oljčnih olj iz 9 italijanskih pokrajin.



Vsako olje je opisano z osmimi podatki, ki jih PCA preslika v 2d ravnino PC_1/PC_2 .



Analiza neznanega olja

6

