

(19)

How many components to keep?
+ How to address the ambiguity of linear recombination?

Of course it depends on goal - denoising, data compression, visualization, modeling

Taking PCA as an example to show basic strategies

Recap

$$Y^{fit} = XA$$

$$Y: m \times r$$

$$X: m \times n$$

$$A: n \times r$$

$$Y^{fit} = \underbrace{U^T}_{"X"} \underbrace{\Lambda}_{"A"} V$$

$$U: n \times m$$

$$\Lambda: n \times n \text{ diag}$$

$$V: n \times r$$

X: "sources"

A: "mixing matrix"

X's columns = eigenvectors of YY^T ,
orthonormal

$A = X^T Y$, rows of A are
(left) eigenvectors of $Y^T Y$,
orthogonal but not normalized

rows of U are left eigenvectors
of $Y Y^T$, orthonormal

rows of V are left eigenvectors of
 $X^T Y$, orthonormal

$$\Lambda = \begin{pmatrix} \sqrt{\lambda_1} & & \\ & \ddots & \\ & & \sqrt{\lambda_n} \end{pmatrix}$$

where λ 's are descending eivals
of YY^T or $Y^T Y$.

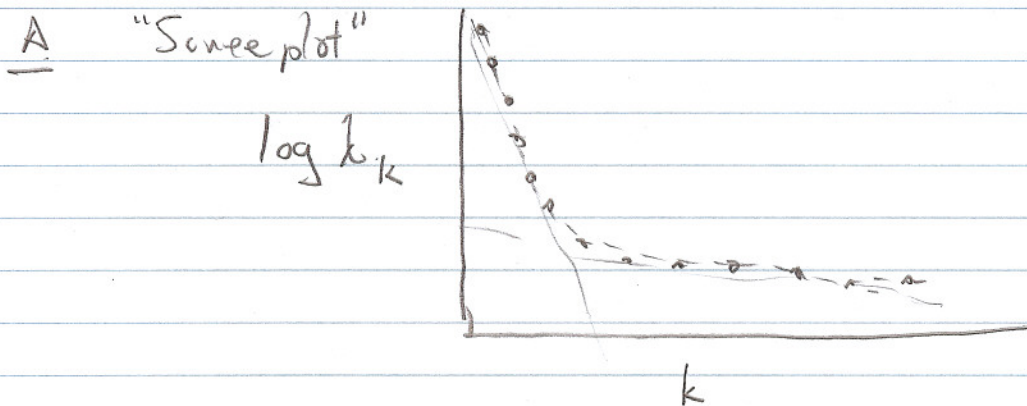
$$\begin{matrix} \downarrow & \downarrow \\ m \times m & r \times r \end{matrix}$$

If $n = \min(m, r)$, then $Y^{fit} = Y$.

Typically, if $n < \min(m, r)$, then some variance is unexplained,
namely, $\sum |Y^{fit} - Y|^2 = \sum \lambda_k$
 $n+1$ to $\min(m, r)$

(20)

Several strategies to choose a smaller value of n ($n = \min(m, r)$):



guaranteed to descend (λ 's in order of importance); local slope indicates fraction of remaining variance explained by next component.

So choose the knee. (arbitrary).

OR, make such plots for random Y 's, & choose the λ 's that are higher than expected for random data.

[Uncorrelated data, \rightarrow expect λ 's all equal for vast amounts of data, (∞)
but finite data will give a range of λ 's.]

[Can also generate surrogate correlated data based on some model of noise correlations]

[Eigenvalues of matrices of correlated Gaussians - a classic problem, many exact results but "difficult"]

B For data reduction/compression - choose some (arbitrary) fraction of variance to explain, (eg, 99%) - keep eigenvs up to that level

(21)

C. (Imaging analysis)

You may have some design variables, eg, S_1, \dots, S_D , each as a time series S_{dim} - and you want to "keep"

the components that are correlated with S_j , i.e., factor

keep λ_h if $\left| \sum_{j=1}^M S_{d,j} X_{j,h} \right| > \text{criterion}$ The retained λ 's might not be consecutive.

\sim Fisher Disc, QIFA

D. (Visualization)

keep 2 or 3 λ 's

The entities $Y^{fit} = XA$ is just as good a fit as

$$Y^{fit} = (XQ)(Q^T A)$$

$$Y^{fit} = U^T \Lambda V \text{ with } UU^T = I_{n \times n}, VV^T = I_{n \times n} \text{ describes}$$

but only in terms of contribution to the variance.

Start with $Y^{fit} = XA$ and "move" normality from X to A :

$$Y^{fit} = (X\Lambda)(\Lambda^{-1}A) \text{, now } A' = \Lambda^{-1}A \text{ has orthogonal normalized rows.}$$

A'_{jk} (or A_{jk}) indicates how much the j^{th} factor (source) (column of X) contributes to the k^{th} variable (column of Y)

Typically, the mixing matrix A' will have most elements non-zero, i.e., most factors contribute to most observed variables

(22)

The "factor rotation" strategy is to find a rotation matrix

R , for which RA' has been "sparsified", i.e.,

by most elements either near 0, or near ± 1 .

Correspondence in Matlab: The TRANSPOSE of the above A' is the

"loadings matrix" A , referred to in Matlab. $A_{\text{MATLAB}} = (A'_{JV})^T$.

VARI-MAX criterion is to maximize (for RA' composed of elements a_{jk})

$$\sum_j \left(\frac{1}{r} \sum_k a_{jk}^4 - \gamma \left(\frac{1}{r} \sum_k a_{jk}^2 \right)^2 \right)$$

$$\gamma = 1 \quad - \text{VARI-MAX}$$

$$\gamma = 0 \quad = \text{"QUARTIMAX"}$$

Since A' is often cols $\frac{1}{r} A'$ + \therefore has normalized rows, so γ -term is constant, independent of rotation.

The rationale here is that "simple" (sparse) mixing matrices are preferable - by making the mixing matrix as far from uniform as possible (i.e., as dispersed as possible)

23

ICA (Independent Components Analysis)

- another approach to the mixing problem

- make the X_i 's as dispersed as possible

"Cocktail Party Problem"

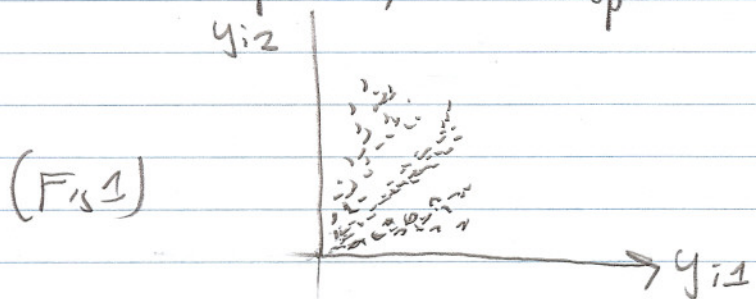
Each speaker $(1, \dots, n)$ creates a time series $(1, \dots, m)$, $X_i (m \times n)$

Each microphone $(1, \dots, r)$ picks up sound from each speaker \hat{c} efficacy $A (n \times r)$.

Can we infer $X + A$ from $Y = XA$? ($Y: m \times r$)

Above ambiguity problem would suggest that we can't - but there's another strategy available.

3 speakers, 2 microphones - we can plot Y_{i1} vs Y_{i2}



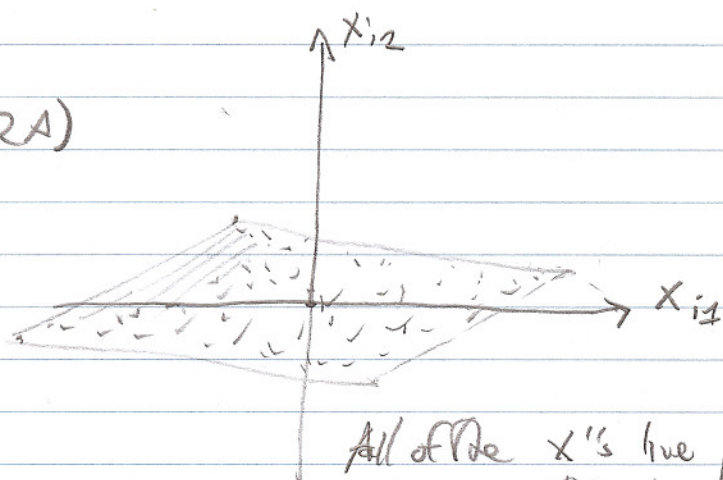
3 "plumes", correspond to how each of the speakers is picked up by each microphone.

PCA will never identify this -- it can't find more than 2 components --

Another motivational example:

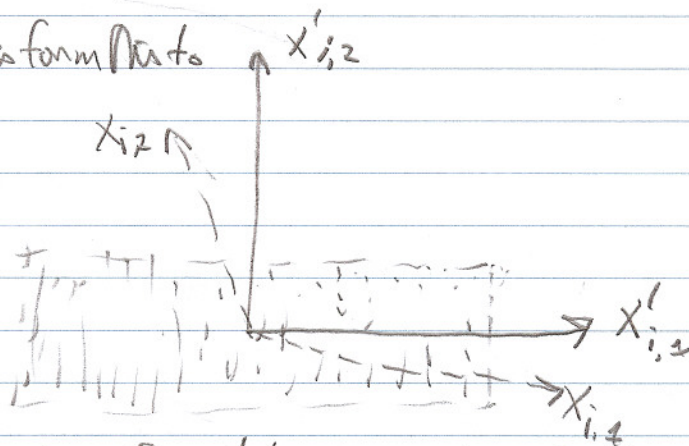
Say Y is $m \times r$, $r \gg 2$, and you find an excellent fit of $Y_{fit} = XA$ for $X = m \times 2$.

(Fig 2A)



All of the x 's live in this parallelogram. $x_{i,1}$ explains most of the variance.

(Fig 2B)



In the new coords, the x 's are independently distributed.

Common denominator in Fig 1, Fig 2 is that we made use of the observation that the points are not distributed in a bunched fashion.

(25)

The idea behind ICA is to find a W for which

WY is "non-Gaussian" the observed Y into components that are as non-Gaussian as possible.

See web tutorial by Aapo Hyvärinen

(Google ICA Tutorial)

All implementations are iterative (like VARIMAX)

One approach: In a Gaussian, if variance is V , then kurtosis

$$K = \left(\frac{\langle (x_j - \langle x_j \rangle)^4 \rangle - 3V^2}{3V^2} \right) \text{ is } 0.$$

So maximize K^2 .

$x_j =$ typical value of j^{th} projection

Another Minimize the entropy of the projection X

Entropy of projection onto X_j is

$$- \sum p_j(x) \log p_j(x)$$

Since a Gaussian has the maximum entropy possible, given a constraint on variance, minimizing the entropy (= maximizing the negentropy) is a reasonable def. of "non-Gaussian".

(26)

Difficult (ie, need lots of data) to estimate entropy from data, based on empirical $p_j(x)$

So - use a parametric approx to $p_j(x)$; this leads back to a moment criterion.

Third approach Minimize the mutual information between the X 's. ("Independent components")

(See $F_1, 2A, B$).

Since total info in multivariable X is constant, minimizing the mutual info of the X 's is equivalent to maximizing their individual infos (= neg entropy) \rightarrow [second approach]

Essentially, ICA looks for non-Gaussian components.

So it will work really well for artifacts of some kind -
eye blinks
EKG

- it is guaranteed to fail for Gaussian (or near-Gaussian) components.

Other general remarks:

- Bridged components may not be independent
- None of the above make real use of dynamics -- same outcome if one time-shuffles the data
(? Formin analysis first
(? "Hierarchical Decomposition" Repucci et al)