

Outline

- ① PCA
- ② Canonical Correlation Analysis

Principal Component Analysis

Sometimes, we require $\|\mathbf{a}_1\| = 1$ and $\langle \mathbf{a}_i, \mathbf{a}_j \rangle = 0$. Thus the problem is to find an interesting set of (orthogonal) **direction vectors** $\{\mathbf{a}_i : i = 1, \dots, p\}$, where the projection scores of \mathbf{X} onto \mathbf{a}_i are useful.

Principal Component Analysis (PCA) is a linear dimension reduction technique that gives a set of direction vectors of maximal **(projected) variances**.

Take $d = 1$. PCA for the distribution of \mathbf{X} finds \mathbf{a}_1 such that

$$\mathbf{a}_1 = \operatorname{argmax}_{\mathbf{a} \in \mathbb{R}^p, \|\mathbf{a}\|=1} \operatorname{Var}(Z_1(\mathbf{a})) \left(= \operatorname{argmax}_{\mathbf{a} \in \mathbb{R}^p, \|\mathbf{a}\|=1} \mathbf{a}' \operatorname{Var}(\mathbf{X}) \mathbf{a} \right),$$

where $Z_1(\mathbf{a}) = a_1 X_1 + \dots + a_p X_p = \mathbf{a}' \mathbf{X}$.

The next section would be

① PCA

② Canonical Correlation Analysis

Linear dimension reduction

For a random vector $\mathbf{X} \in \mathbb{R}^p$, consider reducing the dimension from p to d , i.e., p variables $(X_1, \dots, X_p)^T$ to a set of *most interesting* d variables. Here, $1 \leq d \leq p$.

- Best subset?
- Linear dimension reduction: Construct d variables Z_1, \dots, Z_d as linear combinations of X_1, \dots, X_p , i.e.

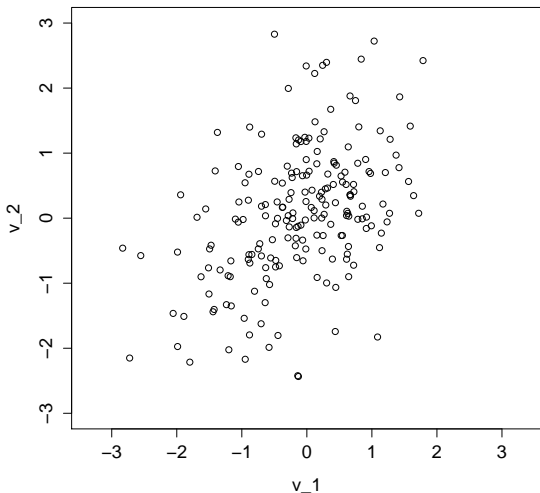
$$Z_i = a_{i1}X_1 + \dots + a_{ip}X_p = \mathbf{a}_i' \mathbf{X} \quad (i = 1, \dots, d),$$

with $\mathbf{a}_i \in \mathbb{R}^p$.

Linear dimension reduction seeks a sequence of such Z_i , or equivalently a sequence of \mathbf{a}_i , where the random variables (Z_i s) are most **important** among all other choices.

Geometric understanding of PCA for point cloud

n (= 200) data points in 2D

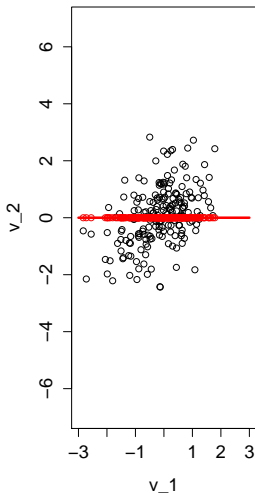


PCA is best understood with a point cloud. Take a look at 2D example.

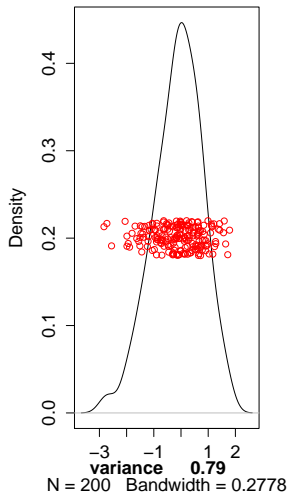
Geometric understanding of PCA for point cloud

Take $\mathbf{a} = (1, 0)'$.

n (= 200) data points in 2D



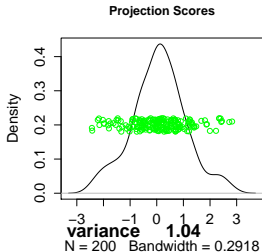
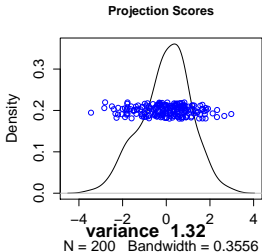
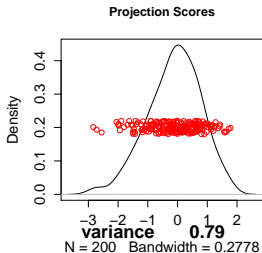
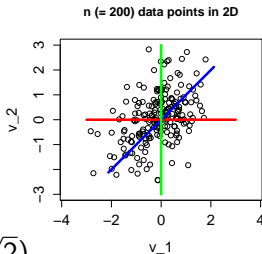
Projection Scores



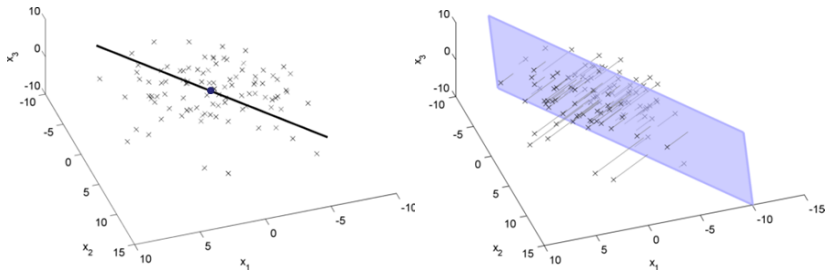
Geometric understanding of PCA for point

Take $\mathbf{a} = (1, 0)'$,
 $(0, 1)'$, $(1/\sqrt{2}, 1/\sqrt{2})'$.

Which
one is better?



Geometric understanding of PCA for 3D point cloud



A 3D point cloud. Mean. 1st PC direction (maximizing variance of projections) explains the cloud best. 1st and 2nd directions form a plane.

Formulation of population PCA-1

Suppose a random vector \mathbf{X} with mean μ , covariance Σ (not necessarily normal).

The first principal component (PC) direction vector is the unit vector $\mathbf{u}_1 \in \mathbb{R}^p$ that maximizes the variance of $\mathbf{u}'_1 \mathbf{X}$ when compared to other unit vectors, i.e.,

$$\mathbf{u}_1 = \operatorname{argmax}_{\mathbf{u} \in \mathbb{R}^p, \|\mathbf{u}\|=1} \operatorname{Var}(\mathbf{u}' \mathbf{X}).$$

- $\mathbf{u}_1 = (u_{11}, \dots, u_{1p})'$ is the first PC direction vector, sometimes called *loading vector*.
- (u_{11}, \dots, u_{1p}) are loadings of the 1st PC.
- $Z_1 = u_{11}X_1 + \dots + u_{1p}X_p = \mathbf{u}'_1 \mathbf{X}$ is the first PC score or the first principal component (random variable).
- $\lambda_1 = \operatorname{Var}(\mathbf{u}' \mathbf{X}) = \operatorname{Var}(Z_1)$ is the variance explained by the first PC.

Formulation of population PCA-2

The second PC direction is the unit vector $\mathbf{u}_2 \in \mathbb{R}^p$ that

- maximizes the variance of $\mathbf{u}_2' \mathbf{X}$;
- is orthogonal to the first PC direction \mathbf{u}_1 .

That is,

$$\mathbf{u}_2 = \underset{\mathbf{u} \in \mathbb{R}^p, \|\mathbf{u}\|=1, \mathbf{u}'\mathbf{u}_1=0}{\operatorname{argmax}} \operatorname{Var}(\mathbf{u}'\mathbf{X}).$$

- $\mathbf{u}_2 = (u_{21}, \dots, u_{2p})'$ is the second PC direction vector, and is the vector of the 2nd set of loadings.
- $Z_2 = \mathbf{u}_2' \mathbf{X}$ is the second principal component.
- $\lambda_2 = \operatorname{Var}(Z_2)$ is the variance explained by the second PC, and $\lambda_1 \geq \lambda_2$.
- $\operatorname{Corr}(Z_1, Z_2) = 0$.

Formulation of population PCA–(3,4,...p)

Given the first $k - 1$ PC directions $\mathbf{u}_1, \dots, \mathbf{u}_{k-1}$, the k th PC direction is the unit vector $\mathbf{u}_k \in \mathbb{R}^p$ that

- maximizes the variance of $\mathbf{u}_k' \mathbf{X}$;
- is orthogonal to the 1st to $(k - 1)$ th PC directions \mathbf{u}_j .

That is,

$$\mathbf{u}_k = \underset{\substack{\mathbf{u} \in \mathbb{R}^p, \|\mathbf{u}\|=1 \\ \mathbf{u}'\mathbf{u}_j=0, j=1, \dots, k-1}}{\operatorname{argmax}} \operatorname{Var}(\mathbf{u}'\mathbf{X}).$$

- $\mathbf{u}_k = (u_{k1}, \dots, u_{kp})'$ is the k th PC direction vector, and is the vector of the k th loadings.
- $Z_k = \mathbf{u}_k' \mathbf{X}$ is the k th principal component.
- $\lambda_k = \operatorname{Var}(Z_k)$ is the variance explained by the k PC score, and $\lambda_1 \geq \dots \geq \lambda_{k-1} \geq \lambda_k$.
- $\operatorname{Corr}(Z_i, Z_j) = 0$ for all $i \neq j \leq k$.

Relation to eigen-decomposition of Σ

Recall the eigen-decomposition of the symmetric positive definite $\Sigma = \mathbf{U}\mathbf{\Lambda}\mathbf{U}'$ with

- $\mathbf{U} = [\mathbf{u}_1, \dots, \mathbf{u}_p]$ orthogonal,
- $\mathbf{\Lambda} = \text{diag}(\lambda_1, \dots, \lambda_p)$ with $\lambda_1 \geq \dots \geq \lambda_p$,
- $\Sigma \mathbf{u}_i = \lambda_i \mathbf{u}_i$.

In next two slides we show that:

- 1 The k th eigenvector \mathbf{u}_k is the k th PC direction vector.
- 2 The k th eigenvalue λ_k is the variance explained by the k th principal component.
- 3 PC directions are both orthogonal $\mathbf{u}_i' \mathbf{u}_j = 0$ ($i \neq j$) and Σ -orthogonal

$$\mathbf{u}_i' \Sigma \mathbf{u}_j = 0 \iff \text{Cov}(Z_i, Z_j) = 0 \quad (i \neq j).$$

Relation to eigen-decomposition of Σ

The first PC direction problem is to maximize $\text{Var}(\mathbf{u}'\mathbf{X})$ with the constraint $\mathbf{u}'\mathbf{u} = 1$. Using Lagrange multiplier λ , the problem of maximization is the same as finding a stationary point of

$$\begin{aligned}\Phi(\mathbf{u}, \lambda) &= \text{Var}(\mathbf{u}'\mathbf{X}) - \lambda(\mathbf{u}'\mathbf{u} - 1) \\ &= \mathbf{u}'\Sigma\mathbf{u} - \lambda(\mathbf{u}'\mathbf{u} - 1).\end{aligned}$$

The stationary point solves the following:

$$\frac{1}{2} \frac{\partial}{\partial \mathbf{u}} \Phi(\mathbf{u}, \lambda) = \Sigma\mathbf{u} - \lambda\mathbf{u} = \mathbf{0},$$

which leads to

$$\lambda = \mathbf{u}'\Sigma\mathbf{u} = \text{Var}(\mathbf{u}'\mathbf{X}), \quad \Sigma\mathbf{u} = \lambda\mathbf{u}. \quad (1)$$

Any eigenvector-eigenvalue pair $(\mathbf{u}_i, \lambda_i)$, $(i = 1, \dots, p)$ satisfies the second eq. in (1). It is clear that the first PC direction is the first eigenvector \mathbf{u}_1 , as it gives the largest variance $\lambda_1 = \mathbf{u}_1'\Sigma\mathbf{u}_1 = \text{Var}(\mathbf{u}_1'\mathbf{X}) \geq \lambda_j$ ($j > 1$).

Relation to eigen-decomposition of Σ

For the k th PC direction, we form a Lagrangian function

$$\Phi(\mathbf{u}, \lambda, \gamma_1^k) = \mathbf{u}'\Sigma\mathbf{u} - \lambda(\mathbf{u}'\mathbf{u} - 1) - \sum_{j=1}^{k-1} 2\gamma_j \mathbf{u}'_j \mathbf{u},$$

given the first $k - 1$ PC directions. The derivative of Φ , equated to zero, is then

$$\begin{aligned} \frac{1}{2} \frac{\partial}{\partial \mathbf{u}} \Phi(\mathbf{u}, \lambda, \gamma_1^k) &= \Sigma\mathbf{u} - \lambda\mathbf{u} - \sum_{j=1}^{k-1} \gamma_j \mathbf{u}_j = \mathbf{0}, \\ \frac{\partial}{\partial \gamma_j} \Phi(\mathbf{u}, \lambda, \gamma_1^k) &= \mathbf{u}'_j \mathbf{u} = 0. \end{aligned} \quad (2)$$

We have $\gamma_j = \mathbf{u}'_j \Sigma \mathbf{u} = 0$ (since $\Sigma \mathbf{u}_j = \lambda_j \mathbf{u}_j$), thus

$$\lambda = \mathbf{u}'\Sigma\mathbf{u} = \text{Var}(\mathbf{u}'\mathbf{X}), \quad \Sigma\mathbf{u} = \lambda\mathbf{u}. \quad (3)$$

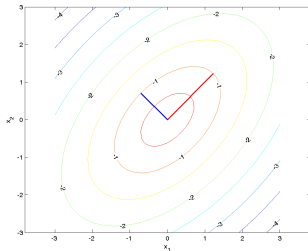
The k th to the last eigen-pairs $(\mathbf{u}_i, \lambda_i)$, $(i = k, \dots, p)$ all satisfy both (3) and (2). Thus, the k th PC direction is \mathbf{u}_k , as it gives the largest variance $\lambda_k = \mathbf{u}'_k \Sigma \mathbf{u}_k$.

Example: Principal components of bivariate normal distribution

Consider $N_2(\mathbf{0}, \Sigma)$ with

$$\Sigma = \begin{bmatrix} 2 & 1 \\ 1 & 2 \end{bmatrix} = \begin{bmatrix} \cos(\pi/4) & -\sin(\pi/4) \\ \sin(\pi/4) & \cos(\pi/4) \end{bmatrix} \begin{bmatrix} 3 & 0 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} \cos(\pi/4) & -\sin(\pi/4) \\ \sin(\pi/4) & \cos(\pi/4) \end{bmatrix}',$$

The PC directions \mathbf{u}_i are the principal axes of ellipsoids, representing the density of MVN, with lengths given by $\sqrt{\lambda_i}$.



Sample PCA

Given multivariate data $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_n]_{p \times n}$, the *sample PCA* sequentially finds orthogonal directions of maximal (projected) sample variance.

- 1 For $\tilde{\mathbf{x}}_i = \mathbf{x}_i - \bar{\mathbf{x}}$, the centered data matrix is

$$\tilde{\mathbf{X}} = [\tilde{\mathbf{x}}_1, \dots, \tilde{\mathbf{x}}_n] = \mathbf{X}(\mathbb{I}_n - \frac{1}{n}\mathbf{1}\mathbf{1}').$$

- 2 The sample variance matrix is then

$$\mathbf{S} = \frac{1}{n-1} \tilde{\mathbf{X}}\tilde{\mathbf{X}}'.$$

- 3 Eigen-decomposition of $\mathbf{S} = \hat{\mathbf{U}}\hat{\mathbf{\Lambda}}\hat{\mathbf{U}}'$ leads to the k th sample PC direction ($\hat{\mathbf{u}}_k$) and the variance of the k th sample score ($\hat{\lambda}_k$).

Sample PCA

It can be checked that the eigenvector $\hat{\mathbf{u}}_k$ of \mathbf{S} satisfies

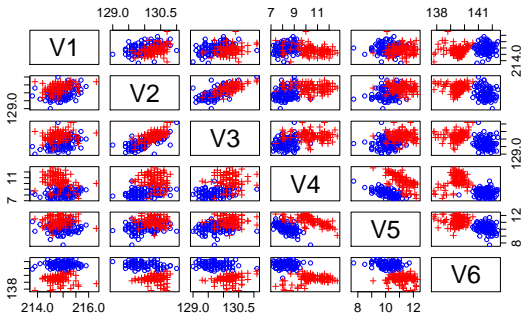
$$\hat{\mathbf{u}}_k = \underset{\substack{\mathbf{u} \in \mathbb{R}^p, \|\mathbf{u}\|=1 \\ \mathbf{u}'\hat{\mathbf{u}}_j=0, j=1, \dots, k-1}}{\operatorname{argmax}} \quad \mathbf{u}'\mathbf{S}\mathbf{u}$$

- $\hat{\mathbf{u}}_k = (\hat{u}_{k1}, \dots, \hat{u}_{kp})'$ is the k th sample PC direction vector, and is the vector of the k th loadings.
- $\mathbf{z}_{(k)} = (z_{(k)1}, \dots, z_{(k)n})_{1 \times n} = \hat{\mathbf{u}}_k' \tilde{\mathbf{X}}$ is the k th score vector, where $z_{(k)i} = \hat{\mathbf{u}}_k' \tilde{\mathbf{x}}_i$.
- $\lambda_k =$ the sample variance of $\{z_{(k)i} : i = 1, \dots, n\}$.
- $\mathbf{u}'\mathbf{S}\mathbf{u} = (n-1)^{-1}(\mathbf{u}'\tilde{\mathbf{X}})(\tilde{\mathbf{X}}'\mathbf{u}) = (n-1)^{-1} \sum_{i=1}^n (\mathbf{u}'\tilde{\mathbf{x}}_i)^2$, which is the sample variance of $\{\mathbf{u}'\mathbf{x}_i, i = 1, \dots, n\}$

The information of the first two principal components is all contained in the score vectors $(\mathbf{z}_{(1)}, \mathbf{z}_{(2)})$, the 2D scatters of which is thus most informative with the largest total variance.

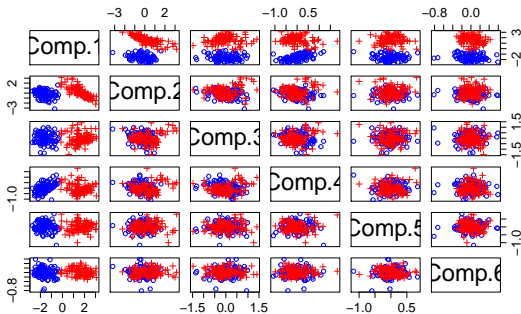
Example: Swiss Bank Note Data

Swiss Bank Note data from HS. $p = 6$, $n = 200$. Genuine (blue) and fake (red) samples in original measurements.



Example: Swiss Bank Note Data

Scatters of principal components



How many components to keep? (1)

The criterion for PCA is a high variance in the principal components. The question involves “how much the PCs explain the variation within the whole data.”

- 1 Total variance in \mathbf{X} is the sum of all marginal variances

$$\begin{aligned}\sum_{k=1}^p \text{Var}(\{x_{ki} : i = 1, \dots, n\}) &= \text{Trace}(\Sigma) \\ &= \sum_{k=1}^p \text{Var}(\{z_{(k)i} : i = 1, \dots, n\}) = \sum_{k=1}^p \lambda_k.\end{aligned}$$

- 2 Variance of the k th scores: $\text{Var}(\{z_{(k)i} : i = 1, \dots, n\}) = \lambda_k$.
- 3 Total variance in the 1st- k th PCs: $\lambda_1 + \dots + \lambda_k$.

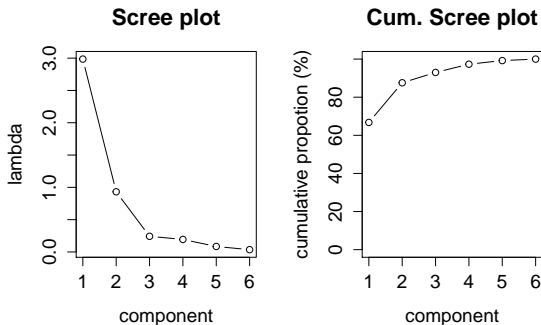
For heuristics made from these quantities, see next slide:

Example: Swiss Bank Note Data–scree plot

In scree plot (k, λ_k) , we look for an elbow.

In cumulative scree plot, (proportion of variance explained,

$(k, \frac{\sum_{j=1}^k \lambda_j}{\sum_{j=1}^p \lambda_j})$), use 90% as a cutoff.



How many components to keep? (2)

- 1 Kaiser's rule (of thumb): Retain PCs 1– k satisfying

$$\lambda_k > \bar{\lambda} = \frac{1}{p} \sum_{j=1}^p \lambda_j.$$

Tends to choose fewer components.

- 2 Likelihood ratio testing on null hypothesis

$$H_0(k) : \lambda_{k+1} = \cdots = \lambda_p,$$

where k components are used if $H_0(k)$ is not rejected at a specified level.

Which variables are most responsible for the principal components?

- Loadings of principal component direction.
- **Biplot** - scatterplot of PC1 and PC2 scores, overlaid with p vectors each representing the loadings of the first two PC directions.

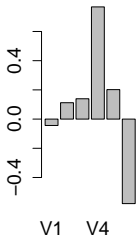
In the Swiss Bank Note Data, the loadings are

	Comp.1	Comp.2	Comp.3	Comp.4	Comp.5	Comp.6
V1			-0.326	0.562	0.753	
V2	0.112		-0.259	0.455	-0.347	-0.767
V3	0.139		-0.345	0.415	-0.535	0.632
V4	0.768	-0.563	-0.218	-0.186		
V5	0.202	0.659	-0.557	-0.451	0.102	
V6	-0.579	-0.489	-0.592	-0.258		

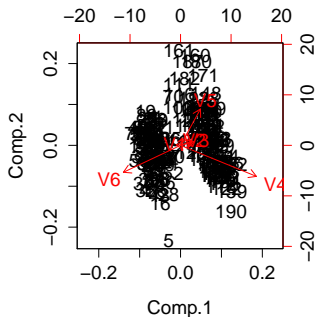
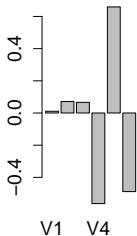
Example: Swiss Bank Note Data–biplot

Which measurements are most responsible for the first two principal components?

PC1 loadings



PC2 loadings



Computation of PCA

PCA is either computed using *eigenvalue decomposition* of $\mathbf{S} = \frac{1}{n-1} \tilde{\mathbf{X}} \tilde{\mathbf{X}}'$ or using the *singular value decomposition* of $\tilde{\mathbf{X}}$.

Eigen-decomposition of \mathbf{S}

For $\mathbf{S} = \mathbf{U} \mathbf{\Lambda} \mathbf{U}'$,

- 1 PC directions \mathbf{u}_k (eigenvectors)
- 2 Variance of PC (scores) λ_k (eigenvalues)
- 3 Matrix of principal component scores

$$\mathbf{V} = \mathbf{U}' \mathbf{X} = \begin{bmatrix} \mathbf{z}_{(1)} \\ \vdots \\ \mathbf{z}_{(p)} \end{bmatrix}.$$

Computation of PCA

Singular value decomposition (SVD) of $\tilde{\mathbf{X}}$

The singular value decomposition (SVD) of $p \times n$ matrix $\tilde{\mathbf{X}}$ has the form

$$\tilde{\mathbf{X}} = \mathbf{U} \mathbf{D} \mathbf{V}'.$$

- The left singular vectors $\mathbf{U} = [\mathbf{u}_1, \dots, \mathbf{u}_p]_{p \times p}$ and the right singular vectors $\mathbf{V} = [\mathbf{v}_1, \dots, \mathbf{v}_p]_{n \times p}$ are orthogonal ($\mathbf{U}'\mathbf{U} = \mathbb{I}_p$, $\mathbf{V}'\mathbf{V} = \mathbb{I}_p$).
- The columns of \mathbf{U} span the column space of $\tilde{\mathbf{X}}$; the columns of \mathbf{V} span the row space.
- $\mathbf{D} = \text{diag}(d_1, \dots, d_p)$, $d_1 \geq d_2 \geq \dots \geq d_p \geq 0$ are the singular values of $\tilde{\mathbf{X}}$.

Computation of PCA

SVD of $\tilde{\mathbf{X}}$

The SVD of $\tilde{\mathbf{X}} = \mathbf{U}\mathbf{D}\mathbf{V}'$. Then

$$\mathbf{S} = \frac{1}{n-1} \tilde{\mathbf{X}}\tilde{\mathbf{X}}' = \frac{1}{n-1} \mathbf{U}\mathbf{D}\mathbf{V}'\mathbf{V}\mathbf{D}\mathbf{U}' = \mathbf{U}\text{diag}\left(\frac{1}{n-1}d_j^2\right)\mathbf{U}'$$

- ① PC directions \mathbf{u}_k (left singular vectors)
- ② Variance of PC (scores) $\frac{1}{n-1}d_j^2$ (singular values²)
- ③ Matrix of principal component scores (right singular vectors)

$$\begin{bmatrix} \mathbf{z}_{(1)} \\ \vdots \\ \mathbf{z}_{(p)} \end{bmatrix} = \mathbf{Z} = \mathbf{U}'\tilde{\mathbf{X}} = \mathbf{D}\mathbf{V}' = \begin{bmatrix} d_1\mathbf{v}'_1 \\ \vdots \\ d_p\mathbf{v}'_p \end{bmatrix}$$

NOTE: we are working with the centered $\tilde{\mathbf{X}}$ here, not \mathbf{X} !!

PCA in R

The standard data format is the $n \times p$ data frame or matrix x .
To perform PCA by eigen decomposition:

```
spr <- princomp(x)
U <- spr$loadings
L <- (spr$sdev)^2
Z <- spr$scores
```

To perform PCA by singular value decomposition

```
gpr <- prcomp(x)
U <- gpr$rotation
L <- (gpr$sdev)^2
Z <- gpr$x
```

Correlation PCA

- PCA is not scale invariant.
- Correlation matrix of a random vector \mathbf{X} is given by

$$\mathbf{R} = D_{\Sigma}^{-\frac{1}{2}} \Sigma D_{\Sigma}^{-\frac{1}{2}},$$

where D_{Σ} is the $p \times p$ diagonal matrix consisting of diagonal elements of Σ .

- Correlation PCA: The PC directions and variance of PC scores are obtained by eigen-decomposition of $\mathbf{R} = \mathbf{U}_R \Lambda_R \mathbf{U}_R'$.
- Preferred if measurements are not commensurate (e.g. X_1 = household income, X_2 = years in school).

PCA for Olivetti Faces data

Olivetti Faces data

- Obtained from <http://www.cs.nyu.edu/~roweis/data.html>.
- Grayscale faces 8 bit [0-255], a few (10) images of several (40) different people.
- 400 total images, 64x64 size.
- From the Olivetti database at ATT.



Images as data

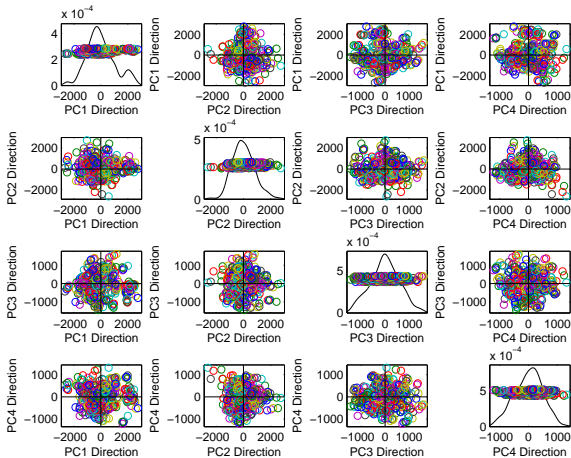
An image is a matrix-valued datum. In Olivetti Faces data, the matrix is of size 64×64 , with each pixel having values between $[0-255]$. The matrix, corresponding one observation, is vectorized (vec'd) by stacking each column into one long vector of size $d = 4096 = 64 \times 64$.

So, \mathbf{x}_1 is a $d \times 1$ vector corresponding to

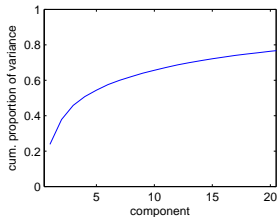
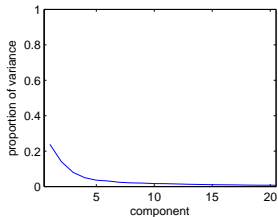
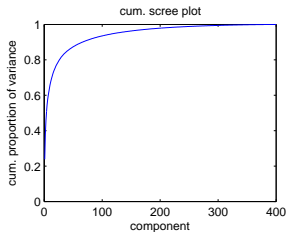
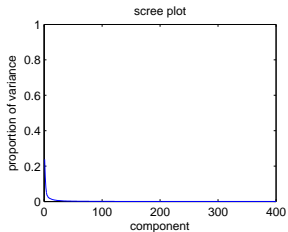


PCA is applied to the data matrix $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_n]$.

Olivetti Faces data—Major components

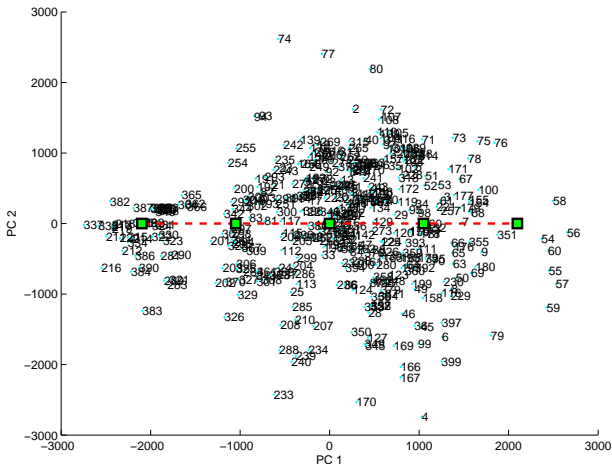


Olivetti Faces data–Scree plots



Olivetti Faces data–Interpretation

Examine the mode of variation by walking along the PC direction from the mean. Here, $\pm 1, 2$ standard deviations of $Z_{(k)}$.



Olivetti Faces data–Interpretation (Eigenfaces)

PC walk. ± 2 std along PC dir. (Top—PC 1, Mid—PC 2, Bottom—PC 3)



PC1 \sim darker to lighter face

PC2 \sim feminin to masculine face

PC3 \sim oval to rectangle face

Olivetti Faces data–Interpretation (Eigenfaces–Computation)

Vectorized data in $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_n]$, with mean $\bar{\mathbf{x}}$. Denote $(\mathbf{u}_j, \lambda_j)$: the j th PC direction and PC variance.

Walk along PC j direction and examine at position $s = \pm 2, \pm 1, 0$ by

- 1 Reconstruction at s : $\mathbf{w}_s = \bar{\mathbf{x}} \pm s\sqrt{\lambda_j}\mathbf{u}_j$
- 2 Convert to image by reshaping the 4096×1 vector \mathbf{w}_s into 64×64 matrix \mathbf{W}_s .

Olivetti Faces data–Interpretation (Eigenfaces–Computation)

Vectorized data in $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_n]$, with mean $\bar{\mathbf{x}}$. Denote $(\mathbf{u}_j, \lambda_j)$: the j th PC direction and PC variance.

Walk along PC j direction and examine at position $s = \pm 2, \pm 1, 0$ by

- 1 Reconstruction at s : $\mathbf{w}_s = \bar{\mathbf{x}} \pm s\sqrt{\lambda_j}\mathbf{u}_j$
- 2 Convert to image by reshaping the 4096×1 vector \mathbf{w}_s into 64×64 matrix \mathbf{W}_s .

Olivetti Faces data–Reconstruction of original data

Approximation to the original data matrix:

$$\mathbf{x}_i = \bar{\mathbf{x}} + \sum_{j=1}^p z_{(j)i} \mathbf{u}_j, \quad (i = 1, \dots, n)$$

Approximation of the original observation \mathbf{x}_i by the first $m < p$ principal components:

$$\hat{\mathbf{x}}_i = \bar{\mathbf{x}} + \sum_{j=1}^m z_{(j)i} \mathbf{u}_j,$$

- The larger m , the better approximation by $\hat{\mathbf{x}}_i$.
- The smaller m , the more succinct dimension reduction of \mathbf{X} .

See some mathematical explanations in the next page.

Olivetti Faces data–Reconstruction of original data

Recall

① $\tilde{\mathbf{X}} = \mathbf{U}\mathbf{D}\mathbf{V}'$

② $\mathbf{Z} = \mathbf{D}\mathbf{V}'$

③ $\tilde{\mathbf{X}} = \mathbf{U}\mathbf{Z}$

④ $\tilde{\mathbf{x}}_i = \mathbf{U}\mathbf{z}_i = \sum_{j=1}^p \mathbf{u}_j z_{(j)i}$

In a coordinate system with $\{\mathbf{u}_i, i = 1, \dots, p\}$ as the p basis vectors, $z_{(j)i}$ is the j th coordinate for the i th observation

$$\tilde{\mathbf{x}}_i = \mathbf{x}_i - \bar{\mathbf{x}}.$$

Reconstruction of original face

Observation index $i = 5$.

5th face. from top left to bottom right: (mean, 1, 5, 10) & (20, 50, 100, 400) PCs



Reconstruction of original face

Observation index $i = 19$.

19th face. from top left to bottom right: (mean, 1, 5, 10) & (20, 50, 100, 400) PCs



Reconstruction of original face

Observation index $i = 100$.

100th face. from top left to bottom right: (mean, 1, 5, 10) & (20, 50, 100, 400) PCs



Olivetti Faces data

Reconstruction of original face

- Human eyes require > 50 principal components to see resemblance between $\hat{\mathbf{x}}_i$ and \mathbf{x}_i .
- Corresponds to about 90 percent of variance explained in PCs.
- Subjective and heuristic decision on “how many components to use”
- Reconstruction by PCA most useful when
 - each datum is visually represented (rather than being just numbers).
 - for example: data objects are images, functions, shapes.

PCA as a mean of dimension reduction

- We can use only the first d PCs to approximately represent the data. Instead of $\mathbf{X}_{p \times n}$, we store the data as $\mathbf{Z}_{d \times n}$.
- This is essentially a dimension reduction approach. However,
 - ① Unsupervised learning (no information on Y).
 - ② Each PC (new variable) is a linear combination of p variables. Hard to interpret.
 - ③ Eigen-decomposition/SVD are problematic when $p \gg n$.

Some open questions

- PCA consistency
 - ① p fixed $n \rightarrow \infty$
 - ② n fixed $p \rightarrow \infty$
 - ③ n and $p \rightarrow \infty$
- Sparse PCA
 - ① Many zeros in \mathbf{u}_k

The next section would be

① PCA

② Canonical Correlation Analysis

Dimension reduction for two sets of random variables

When distinction between explanatory and response variables are not so clear, an analysis dealing with the two sets of variables in a symmetric manner is desired.

Examples

- ① Relationship between genes expressions and biological variables;
- ② Relation between two sets of psychological tests, each with multidimensional measurements.

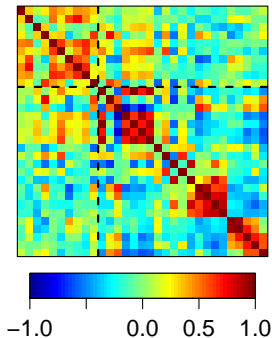
Example: nutrimouse dataset

- a nutrition study in the mouse
- Reference: <http://www6.toulouse.inra.fr/toxalim/pharmaco-moleculaire/acceuil.html>
Pascal Martin from the Toxicology and Pharmacology Laboratory (French National Institute for Agronomic Research).
- Obtained from R package, CCA.
- $n = 40$ mice, each with $r = 120$ gene expression levels, associated with nutrition problem, and $s = 21$ measurements of concentrations of 21 hepatic fatty acids.

$$\mathbf{X}_{120 \times 40}, \quad \mathbf{Y}_{21 \times 40}$$

Example: nutrimouse dataset–Objective

XY correlation

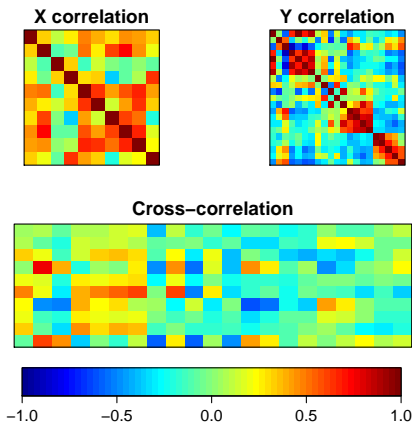


Focus on the first 10 gene expression levels (Dimension r of the first set of variables is now 10).

$$\mathbf{X}_{10 \times 40}, \quad \mathbf{Y}_{21 \times 40}$$

Example: nutrimouse dataset–Objective

- Interested in dimension reduction of \mathbf{X} and \mathbf{Y} , while keeping the important association between \mathbf{X} and \mathbf{Y} .
- Find linear dimension reduction of \mathbf{X} and \mathbf{Y} using the cross-correlation matrix



Canonical Correlation Analysis (CCA)

CCA seeks to identify and quantify the associations between two sets of variables.

Given two random vectors $\mathbf{X} \in \mathbb{R}^r$ and $\mathbf{Y} \in \mathbb{R}^s$ ($r = s$ or $r \neq s$), consider *linear dimension reduction* of each of two random vectors,

$$\begin{aligned}\xi &= \mathbf{g}'\mathbf{X} = g_1X_1 + \cdots + g_rX_r, \\ \omega &= \mathbf{h}'\mathbf{Y} = h_1Y_1 + \cdots + h_sY_s.\end{aligned}$$

CCA finds the random variables (ξ, ω) or the projection vectors (\mathbf{g}, \mathbf{h}) that give *maximal correlation* between ξ and ω ,

$$\text{Corr}(\xi, \omega) = \frac{\text{Cov}(\xi, \omega)}{\sqrt{\text{Var}(\xi)\text{Var}(\omega)}}$$

Population CCA – 1

Denote

$$\begin{aligned}\Sigma_{11} &= \text{Cov}(\mathbf{X}), \Sigma_{22} = \text{Cov}(\mathbf{Y}) \\ \Sigma_{12} &= \text{Cov}(\mathbf{X}, \mathbf{Y}), \Sigma_{21} = \text{Cov}(\mathbf{Y}, \mathbf{X}).\end{aligned}$$

The first set of projection vectors

$$(\mathbf{g}_1, \mathbf{h}_1) = \underset{\mathbf{g} \in \mathbb{R}^r, \mathbf{h} \in \mathbb{R}^s}{\operatorname{argmax}} \operatorname{Corr}(\mathbf{g}'\mathbf{X}, \mathbf{h}'\mathbf{Y}) \quad (4)$$

$$\underset{\mathbf{g} \in \mathbb{R}^r, \mathbf{h} \in \mathbb{R}^s}{\operatorname{argmax}} \frac{\mathbf{g}'\Sigma_{12}\mathbf{h}}{\sqrt{\mathbf{g}'\Sigma_{11}\mathbf{g}\mathbf{h}'\Sigma_{22}\mathbf{h}}}. \quad (5)$$

- $(\mathbf{g}_1, \mathbf{h}_1)$ Canonical correlation vectors.
- $(\xi_1 = \mathbf{g}_1'\mathbf{X}, \omega_1 = \mathbf{h}_1'\mathbf{Y})$ are called canonical variates (like PC).
- $\rho_1 = \operatorname{Corr}(\xi_1, \omega_1)$ is the largest canonical correlation.

Population CCA – 2,3,...

The k th set of projection vectors, given $(\mathbf{g}_1, \dots, \mathbf{g}_{k-1})$ and $(\mathbf{h}_1, \dots, \mathbf{h}_{k-1})$, are

$$(\mathbf{g}_k, \mathbf{h}_k) = \underset{\substack{\mathbf{g} \in \mathbb{R}^r, \mathbf{h} \in \mathbb{R}^s \\ \mathbf{g}'\Sigma_{11}\mathbf{g}_j=0, \\ \mathbf{h}'\Sigma_{22}\mathbf{h}_j=0, j=1, \dots, k-1}}{\operatorname{argmax}} \frac{\mathbf{g}'\Sigma_{12}\mathbf{h}}{\sqrt{\mathbf{g}'\Sigma_{11}\mathbf{g}\mathbf{h}'\Sigma_{22}\mathbf{h}}}. \quad (6)$$

- $(\mathbf{g}_k, \mathbf{h}_k)$ Canonical correlation vectors.
- $(\xi_k = \mathbf{g}_k' \mathbf{X}, \omega_k = \mathbf{h}_k' \mathbf{Y})$ the k th canonical variates.
- $\rho_k = \operatorname{Corr}(\xi_k, \omega_k) \leq \rho_j, j = 1, \dots, k-1$.
- $\operatorname{Corr}(\xi_k, \xi_j) = 0, \operatorname{Corr}(\omega_k, \omega_j) = 0, j = 1, \dots, k-1$.
- Generally $\mathbf{g}_i' \mathbf{g}_j = 0$ NOT true.

Sample CCA

For n sample $(\mathbf{x}_1, \mathbf{y}_i)$, $(i = 1, \dots, n)$, denote

$$\begin{aligned}\mathbf{S}_{11} &= \widehat{\text{Cov}}(\mathbf{X}), \mathbf{S}_{22} = \widehat{\text{Cov}}(\mathbf{Y}) \\ \mathbf{S}_{12} &= \widehat{\text{Cov}}(\mathbf{X}, \mathbf{Y}), \mathbf{S}_{21} = \widehat{\text{Cov}}(\mathbf{Y}, \mathbf{X}).\end{aligned}$$

The sample CCA is

$$(\mathbf{g}_k, \mathbf{h}_k) = \underset{\substack{\mathbf{g} \in \mathbb{R}^r, \mathbf{h} \in \mathbb{R}^s \\ \mathbf{g}'\mathbf{S}_{11}\mathbf{g}_j=0, \\ \mathbf{h}'\mathbf{S}_{22}\mathbf{h}_j=0, j=1, \dots, k-1}}{\text{argmax}} \frac{\mathbf{g}'\mathbf{S}_{12}\mathbf{h}}{\sqrt{\mathbf{g}'\mathbf{S}_{11}\mathbf{g}\mathbf{h}'\mathbf{S}_{22}\mathbf{h}}}. \quad (7)$$

Finding CCA solution

First CC vectors: maximize $\frac{\mathbf{g}'\mathbf{S}_{12}\mathbf{h}}{\sqrt{\mathbf{g}'\mathbf{S}_{11}\mathbf{g}\mathbf{h}'\mathbf{S}_{22}\mathbf{h}}}$.

Change-of-variable: $\mathbf{a} = \mathbf{S}_{11}^{\frac{1}{2}}\mathbf{g}$, $\mathbf{b} = \mathbf{S}_{22}^{\frac{1}{2}}\mathbf{h}$. The problem is now to maximize

$$\frac{\mathbf{a}'\mathbf{S}_{11}^{-\frac{1}{2}}\mathbf{S}_{12}\mathbf{S}_{22}^{-\frac{1}{2}}\mathbf{b}}{\sqrt{\mathbf{a}'\mathbf{a}\mathbf{b}'\mathbf{b}}}$$

with respect to $\mathbf{a} \in \mathbb{R}^r$, $\mathbf{b} \in \mathbb{R}^s$, and the solution is given by

$$\mathbf{a}_1 = \mathbf{v}_1(\mathbf{S}_{11}^{-\frac{1}{2}}\mathbf{S}_{12}\mathbf{S}_{22}^{-1}\mathbf{S}_{21}\mathbf{S}_{11}^{-\frac{1}{2}}),$$

$$\mathbf{b}_1 = \mathbf{v}_1(\mathbf{S}_{22}^{-\frac{1}{2}}\mathbf{S}_{21}\mathbf{S}_{11}^{-1}\mathbf{S}_{12}\mathbf{S}_{22}^{-\frac{1}{2}}),$$

where $\mathbf{v}_i(\mathbf{M})$ is the i th eigenvector (corresponding to the i th largest eigenvalue) of symmetric \mathbf{M} .

Finding CCA solution

For $j = 1, \dots, \min(s, r)$,

$$\mathbf{a}_j = \mathbf{v}_j(\mathbf{S}_{11}^{-\frac{1}{2}} \mathbf{S}_{12} \mathbf{S}_{22}^{-1} \mathbf{S}_{21} \mathbf{S}_{11}^{-\frac{1}{2}}),$$

$$\mathbf{b}_j = \mathbf{v}_j(\mathbf{S}_{22}^{-\frac{1}{2}} \mathbf{S}_{21} \mathbf{S}_{11}^{-1} \mathbf{S}_{12} \mathbf{S}_{22}^{-\frac{1}{2}}),$$

we have

$$\mathbf{g}_j = \mathbf{S}_{11}^{-\frac{1}{2}} \mathbf{a}_j, \quad \mathbf{h}_j = \mathbf{S}_{22}^{-\frac{1}{2}} \mathbf{b}_j$$

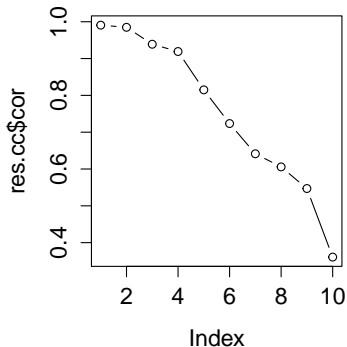
and thus canonical covariate scores $\boldsymbol{\xi}_j = \mathbf{g}_j' \mathbf{X} = (\mathbf{g}_j' \mathbf{x}_1, \dots, \mathbf{g}_j' \mathbf{x}_n)$
and $\boldsymbol{\omega}_j = \mathbf{h}_j' \mathbf{Y} = (\mathbf{h}_j' \mathbf{y}_1, \dots, \mathbf{h}_j' \mathbf{y}_n)$. Note that

- $\mathbf{g}_j' \mathbf{g}_k = \mathbf{a}_j' \mathbf{S}_{11}^{-1} \mathbf{a}_k \neq 0$ (not orthogonal);
- $\widehat{\text{Cov}}(\{\xi_{j(i)}\}, \{\xi_{k(i)}\}) = \mathbf{g}_j' \mathbf{S}_{11} \mathbf{g}_k = 0$ (uncorrelated).

nutrimouse data

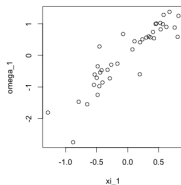
Canonical correlation coefficients

- Reduced dataset: $r = 10$ and $s = 21$ variables with $n = 40$ samples.
- $\min(s, r) = 10$ pairs of canonical variates, with decreasing canonical correlation coefficients.

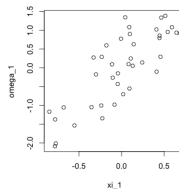


- First four pairs of canonical variates
 $\{(\xi_{(j)i}, \omega_{(j)i}) : i = 1, \dots, n\}$.

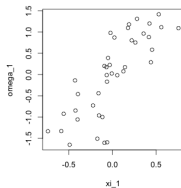
first canonical covariates



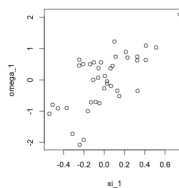
second canonical covariates



third canonical covariates



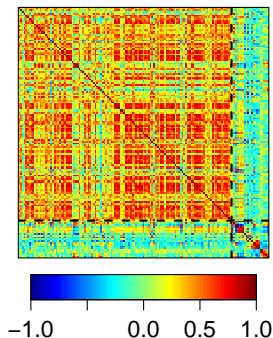
fourth canonical covariates



nutrimouse data with $r = 120$

Original data, with $r \gg n$. Problem?

XY correlation



Regularized CCA

When $r \geq n$, and $s \geq n$, there always exist

$$(\mathbf{g}_1, \mathbf{h}_1) = \operatorname{argmax}_{\mathbf{g} \in \mathbb{R}^r, \mathbf{h} \in \mathbb{R}^s} \frac{\mathbf{g}' \mathbf{S}_{12} \mathbf{h}}{\sqrt{\mathbf{g}' \mathbf{S}_{11} \mathbf{g} \mathbf{h}' \mathbf{S}_{22} \mathbf{h}}},$$

satisfying $\operatorname{Corr}(\xi_j, \omega_k) = \pm 1$ (perfect correlation!). This is an artifact from that \mathbf{S}_{11} and \mathbf{S}_{22} are not invertible (or the rank of $\tilde{\mathbf{X}}$ and $\tilde{\mathbf{Y}}$ is at most $n - 1$.)

A way to circumvent this issue and obtain a meaningful estimate of population canonical correlation coefficients, Regularized CCA is often used:

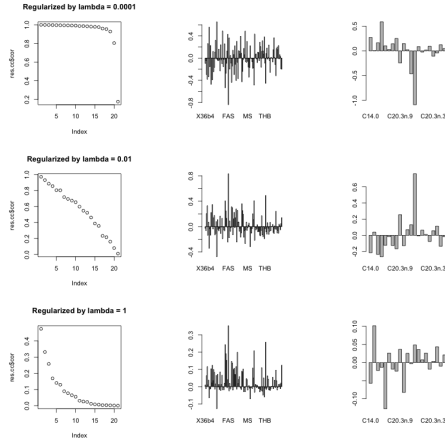
$$(\mathbf{g}_k, \mathbf{h}_k) = \operatorname{argmax}_{\substack{\mathbf{g} \in \mathbb{R}^r, \mathbf{h} \in \mathbb{R}^s \\ \mathbf{g}' \mathbf{S}_{11} \mathbf{g} = 0, \\ \mathbf{h}' \mathbf{S}_{22} \mathbf{h} = 0, j=1, \dots, k-1}} \frac{\mathbf{g}' \mathbf{S}_{12} \mathbf{h}}{\sqrt{\mathbf{g}' (\mathbf{S}_{11} + \lambda_1 \mathbb{I}_r) \mathbf{g} \mathbf{h}' (\mathbf{S}_{22} + \lambda_2 \mathbb{I}_s) \mathbf{h}}}, \quad (8)$$

for $\lambda_1, \lambda_2 > 0$.

nutrimouse data–Regularized CCA

- ① sample canonical correlation coefficients ρ_j and
- ② corresponding projection vectors \mathbf{h}_1 and \mathbf{g}_1

are varying depending on the choice of regularization parameter $\lambda_1 = \lambda_2 = 0.0001, 0.01, 1$.



CCA in R

$n \times r$ data frame or matrix x and $n \times s$ data frame or matrix y
To perform CCA:

```
cc=cancor(x,y)
rho<-cc$cor
g<-cc$xcoef
h<-cc$ycoef
```

To perform regularized CCA, use package CCA,

```
library(CCA)
cc=rcc(x,y,lambd1,lambd2)
rho<-cc$cor
g<-cc$xcoef
h<-cc$ycoef
```

Possible research topics

Sparse PCA

- H Shen, JZ Huang - Journal of multivariate analysis, 2008
- A d'Aspremont, L El Ghaoui, MI Jordan, GRG Lanckriet - SIAM review, 2007
- Z Ma - The Annals of Statistics, 2013
- TT Cai, Z Ma, Y Wu - The Annals of Statistics, 2013

PCA Consistency

- S Jung, JS Marron - The Annals of Statistics, 2009
- D Shen, H Shen, JS Marron - Journal of Multivariate Analysis, 2013
- S Jung, A Sen, JS Marron - Journal of Multivariate Analysis, 2012