

Sparse PCA

Archer Yang

McGill University

October 3, 2024

- **Reading assignment:** Zou (2006); Zou (2018)

PRECURSOR

Motivation of Sparse PCA

Pitprops Data: $N = 180$ and $p = 13$. SPCA generated very **sparse** loading structures *without losing much variance.*

	PCA			SPCA		
	PC1	PC2	PC3	PC1	PC2	PC3
topdiam	-.404	.218	-.207	-.477		
length	-.406	.186	-.235	-.476		
moist	-.124	.541	.141		.785	
testsg	-.173	.456	.352		.620	
ovensg	-.057	-.170	.481	.177		.640
ringtop	-.284	-.014	.475			.589
ringbut	-.400	-.190	.253	-.250		.492
bowmax	-.294	-.189	-.243	-.344	-.021	
bowdist	-.357	.017	-.208	-.416		
whorls	-.379	-.248	-.119	-.400		
clear	.011	.205	-.070			
knots	.115	.343	.092		.013	
diaknot	.113	.309	-.326			-.015
variance	32.4	18.3	14.4	28.0	14.0	13.3

SCoTLASSO

- Jolliffe, Trendafilov, and Uddin (2003) proposed **SCoTLASSO**, which successively maximizes

$$\begin{aligned} & \max_{\|\mathbf{b}_k\|=1} \mathbf{b}_k^\top \mathbf{S} \mathbf{b}_k \\ \text{s.t. } & \mathbf{b}_j^\top \mathbf{b}_k = 0, \quad h < k, \quad \|\mathbf{b}_k\|_1 \leq t \end{aligned}$$

- Nonconvex optimization problem, high computational cost.

INSIGHT

PCA – equivalent form

- SPCA extends the reconstruction view of the PCA to derive sparse PCs. Recall that the first PC v_1 can be defined as

$$\begin{aligned} v_1 = & \arg \min_b \frac{1}{N} \sum_{n=1}^N \|x_n - bb^\top x_n\|^2 \\ & \text{subject to } \|b\|^2 = 1 \end{aligned} \tag{1}$$

- We reformulate (1) as

$$\begin{aligned} v_1 = & \arg \min_{a,b} \frac{1}{N} \sum_{n=1}^N \|x_n - ab^\top x_n\|^2 \\ & \text{subject to } \|a\|^2 = 1 \text{ and } a = b \end{aligned}$$

PCA – equivalent form

- The following theorem says that we can drop $\mathbf{a} = \mathbf{b}$ and still recover \mathbf{v}_1 exactly.

Theorem 1

For any $\lambda > 0$, let

$$(\hat{\mathbf{a}}, \hat{\mathbf{b}}) = \arg \min_{\mathbf{a}, \mathbf{b}} \sum_{n=1}^N \|\mathbf{x}_n - \mathbf{ab}^\top \mathbf{x}_n\|^2 + \lambda \|\mathbf{b}\|^2$$

subject to $\|\mathbf{a}\|^2 = 1$

Then $\hat{\mathbf{b}} \propto \mathbf{v}_1$

PCA – equivalent form

Theorem 2

Let $\widetilde{\mathbf{B}} = (\mathbf{v}_1, \dots, \mathbf{v}_K)$

$$\begin{aligned}\widetilde{\mathbf{B}} = & \arg \min_{\mathbf{B}} \sum_{n=1}^N \|\mathbf{x}_n - \mathbf{B}\mathbf{B}^\top \mathbf{x}_n\|^2 \\ & \text{subject to } \mathbf{B}^\top \mathbf{B} = \mathbf{I}_k\end{aligned}$$

One can show that, for any $\lambda > 0$, let $\widehat{\mathbf{A}} = (\hat{\mathbf{a}}_1, \dots, \hat{\mathbf{a}}_k) \in \mathbb{R}^{p \times k}$ and
 $\widehat{\mathbf{B}} = (\hat{\mathbf{b}}_1, \dots, \hat{\mathbf{b}}_k)$

$$\begin{aligned}(\widehat{\mathbf{A}}, \widehat{\mathbf{B}}) = & \arg \min_{\mathbf{A}, \mathbf{B}} \sum_{n=1}^N \|\mathbf{x}_n - \mathbf{AB}^\top \mathbf{x}_n\|^2 + \lambda \sum_{j=1}^k \|\mathbf{b}_j\|^2 \\ & \text{subject to } \mathbf{A}^\top \mathbf{A} = \mathbf{I}_k\end{aligned}$$

Then $\hat{\mathbf{a}}_j = \mathbf{v}_j$ and $\hat{\mathbf{b}}_j \propto \mathbf{v}_j$ for $j = 1, 2, \dots, k$.

Proof

- Use the notation $\mathbf{A} = (\mathbf{a}_1, \dots, \mathbf{a}_k) \in \mathbb{R}^{p \times k}$ and $\mathbf{B} = (\mathbf{b}_1, \dots, \mathbf{b}_k) \in \mathbb{R}^{p \times k}$.
- Let \mathbf{A}_\perp be any orthonormal matrix such that $\overline{\mathbf{A}} = [\mathbf{A}, \mathbf{A}_\perp]$ is $p \times p$ square orthonormal.
- For any matrix \mathbf{X} and orthogonal matrices \mathbf{A} and \mathbf{A}_\perp :

$$\begin{aligned}\|\mathbf{X}\mathbf{A}\|_F^2 + \|\mathbf{X}\mathbf{A}_\perp\|_F^2 &= \|\mathbf{X}\overline{\mathbf{A}}\|_F^2 = \text{tr}(\overline{\mathbf{A}}^\top \mathbf{X}^\top \mathbf{X} \overline{\mathbf{A}}) \\ &= \text{tr}(\overline{\mathbf{A}\mathbf{A}}^\top \mathbf{X}^\top \mathbf{X}) = \text{tr}(\mathbf{X}^\top \mathbf{X}) \\ &= \|\mathbf{X}\|_F^2\end{aligned}$$

Proof

- We have

$$\begin{aligned} \sum_{n=1}^N \|x_n - AB^\top x_n\|^2 &= \|X - XBA^\top\|_F^2 \\ &= \text{tr}((X^\top - AB^\top X^\top)(X - XBA^\top)) \\ &= \text{tr}(X^\top X - 2AB^\top X^\top X + AB^\top X^\top XBA^\top) \\ &= \|XA_\perp\|_F^2 + \text{tr}(A^\top X^\top XA - 2AB^\top X^\top X \\ &\quad + BA^\top AB^\top X^\top X) \\ &= \|XA_\perp\|_F^2 + \text{tr}(A^\top X^\top XA - 2AB^\top X^\top X \\ &\quad + B^\top X^\top XB) \\ &= \|XA_\perp\|_F^2 + \|XA - XB\|_F^2 \end{aligned}$$

- $\text{tr}(AB^\top X^\top XBA^\top) = \text{tr}(BA^\top AB^\top X^\top X)$, since AB^\top is a square matrix.

Proof

- Let

$$C_\lambda(\mathbf{A}, \mathbf{B}) = \sum_{n=1}^N \|\mathbf{x}_n - \mathbf{AB}^\top \mathbf{x}_n\|^2 + \lambda \sum_{j=1}^k \|\mathbf{b}_j\|^2$$

- Hence, with \mathbf{A} fixed, solving

$$\arg \min_{\mathbf{B}} C_\lambda(\mathbf{A}, \mathbf{B})$$

is equivalent to solving the series of ridge regressions

$$\arg \min_{\mathbf{B}=(\mathbf{b}_j)_{j=1}^k} \sum_{j=1}^k \|\mathbf{X}\mathbf{a}_j - \mathbf{X}\mathbf{b}_j\|^2 + \lambda \|\mathbf{b}_j\|^2$$

Proof

- It is easy to show that

$$\widehat{\mathbf{B}} = (\mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^\top \mathbf{X} \mathbf{A}$$

- Then using Lemma 1

$$C_\lambda(\mathbf{A}, \widehat{\mathbf{B}}) = \text{tr}((\mathbf{X}\mathbf{A})^\top (\mathbf{I} - \mathbf{S}_\lambda)(\mathbf{X}\mathbf{A}))$$

- Rearranging the terms, we get

$$C_\lambda(\mathbf{A}, \widehat{\mathbf{B}}) = \text{tr}(\mathbf{X}^\top \mathbf{X}) - \text{tr}(\mathbf{A}^\top \mathbf{X}^\top \mathbf{S}_\lambda \mathbf{X} \mathbf{A})$$

which must be minimized with respect to \mathbf{A} with $\mathbf{A}^\top \mathbf{A} = \mathbf{I}$.

- This is equivalent to

$$\min_{\mathbf{A}} C_\lambda(\mathbf{A}, \mathbf{B}) \iff \max_{\mathbf{A}} \text{tr}(\mathbf{A}^\top \mathbf{X}^\top \mathbf{S}_\lambda \mathbf{X} \mathbf{A})$$

Proof

- Hence \mathbf{A} should be taken to be the largest k eigenvectors of $\mathbf{X}^\top \mathbf{S}_\lambda \mathbf{X}$.
- If the SVD of \mathbf{X} is

$$\mathbf{X} = \mathbf{U}\mathbf{D}\mathbf{V}^\top$$

- It is easy to show that $\mathbf{X}^\top \mathbf{S}_\lambda \mathbf{X}$ has eigen-decomposition

$$\mathbf{X}^\top \mathbf{S}_\lambda \mathbf{X} = \mathbf{V} \mathbf{D}^2 (\mathbf{D}^2 + \lambda \mathbf{I})^{-1} \mathbf{D}^2 \mathbf{V}^\top$$

hence

$$\widehat{\mathbf{A}} = \mathbf{V}[:, 1:k] = (\mathbf{v}_1, \dots, \mathbf{v}_k)$$

Proof

- Now plugging $\mathbf{X} = \mathbf{U}\mathbf{D}\mathbf{V}^\top$ into

$$\begin{aligned}\widehat{\mathbf{B}} &= (\mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^\top \mathbf{X} \mathbf{A} \\ &= \mathbf{V} (\mathbf{D}^2 + \lambda \mathbf{I})^{-1} \mathbf{D}^2 \mathbf{V}^\top \mathbf{I}_p[, 1:k]\end{aligned}$$

- Therefore $\widehat{\mathbf{B}} \propto \mathbf{V}[, 1:k]$

Lemma 1

Consider the ridge regression criterion

$$C_\lambda(\beta) = \|\mathbf{y} - \mathbf{X}\beta\|^2 + \lambda\|\beta\|^2$$

Then if $\hat{\beta} = \arg \min_{\beta} C_\lambda(\beta)$

$$C_\lambda(\hat{\beta}) = \mathbf{y}^\top (\mathbf{I} - \mathbf{S}_\lambda) \mathbf{y}$$

where \mathbf{S}_λ is the ridge operator

$$\mathbf{S}_\lambda = \mathbf{X}(\mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^\top$$

See proof in Zou (2006)

SPARSE PCA

Sparse PCA

The sparse PCA criterion for the first k sparse principal components is defined as

$$(\widehat{\mathbf{A}}, \widehat{\mathbf{B}}) = \arg \min_{\mathbf{A}, \mathbf{B}} \sum_{n=1}^N \|\mathbf{x}_n - \mathbf{AB}^\top \mathbf{x}_n\|^2 + \lambda \sum_{j=1}^k \|\mathbf{b}_j\|^2 + \lambda_1 \sum_{j=1}^k \|\mathbf{b}_j\|_1$$

subject to $\mathbf{A}^\top \mathbf{A} = \mathbf{I}_k$

Computation of Sparse PCA

- If A is fixed, the optimization problem of B is

$$B^+ = \arg \min_{A, B} \sum_{n=1}^N \|x_n - AB^\top x_n\|^2 + \lambda \sum_{j=1}^k \|b_j\|^2 + \lambda_1 \sum_{j=1}^k \|b_j\|_1$$

- For $B^+ = (b_1^+, \dots, b_k^+)$, this can be decomposed into

$$b_j^+ = \arg \min_{A, B} \|Xa_j - Xb_j\|^2 + \lambda \sum_{j=1}^k \|b_j\|^2 + \lambda_1 \sum_{j=1}^k \|b_j\|_1$$

- The above is an elastic net penalized least square problem

Computation of Sparse PCA

- If B is fixed, the optimization problem of A is

$$\mathbf{A}^+ = \arg \min_{\mathbf{A}, \mathbf{B}} \sum_{n=1}^N \left\| \mathbf{x}_n - \mathbf{AB}^\top \mathbf{x}_n \right\|^2$$

$$\text{subject to } \mathbf{A}^\top \mathbf{A} = \mathbf{I}_k$$

- We can see that when $k = p$, it is the Procrustes rotation problem.
- Analytical solution: compute the SVD $(\mathbf{X}^\top \mathbf{X})\mathbf{B} = \mathbf{U}\mathbf{D}\mathbf{V}^\top$ and set

$$\mathbf{A}^+ = \mathbf{U}\mathbf{V}^\top$$

Computation of Sparse PCA

- The SPCA algorithm iterates between the elastic net regression step and the SVD step till convergence.
- The output is the normalized B matrix for $j = 1, \dots, k$

$$\mathbf{v}_j = \hat{\mathbf{b}}_j / \|\hat{\mathbf{b}}_j\|$$

Example: Pitprops Data

Table 1. Pitprops Data: Loadings of the First Six Principal Components

Variable	PC1	PC2	PC3	PC4	PC5	PC6
topdiam	-0.404	0.218	-0.207	0.091	-0.083	0.120
length	-0.406	0.186	-0.235	0.103	-0.113	0.163
moist	-0.124	0.541	0.141	-0.078	0.350	-0.276
testsg	-0.173	0.456	0.352	-0.055	0.356	-0.054
ovensg	-0.057	-0.170	0.481	-0.049	0.176	0.626
ringtop	-0.284	-0.014	0.475	0.063	-0.316	0.052
ringbut	-0.400	-0.190	0.253	0.065	-0.215	0.003
bowmax	-0.294	-0.189	-0.243	-0.286	0.185	-0.055
bowdist	-0.357	0.017	-0.208	-0.097	-0.106	0.034
whorls	-0.379	-0.248	-0.119	0.205	0.156	-0.173
clear	0.011	0.205	-0.070	-0.804	-0.343	0.175
knots	0.115	0.343	0.092	0.301	-0.600	-0.170
diaknot	0.113	0.309	-0.326	0.303	0.080	0.626
Variance (%)	32.4	18.3	14.4	8.5	7.0	6.3
Cumulative variance (%)	32.4	50.7	65.1	73.6	80.6	86.9

Example: Pitprops Data

Table 2. Pitprops Data: Loadings of the First Six Modified PCs by SCoTLASS. Empty cells have zero loadings.

$t = 1.75$ <i>Variable</i>	<i>PC1</i>	<i>PC2</i>	<i>PC3</i>	<i>PC4</i>	<i>PC5</i>	<i>PC6</i>
topdiam	0.664			-0.025	0.002	-0.035
length	0.683	-0.001		-0.040	0.001	-0.018
moist		0.641	0.195		0.180	-0.030
testsg		0.701	0.001			-0.001
ovensg					-0.887	-0.056
ringtop		0.293	-0.186		-0.373	0.044
ringbut	0.001	0.107	-0.658		-0.051	0.064
bowmax	0.001			0.735	0.021	-0.168
bowdist	0.283					-0.001
whorls	0.113		-0.001	0.388	-0.017	0.320
clear						-0.923
knots		0.001		-0.554	0.016	0.004
diaknot			0.703	0.001	-0.197	0.080
Number of nonzero loadings	6	6	6	6	10	13
Variance (%)	19.6	16.0	13.1	13.1	9.2	9.0
Adjusted variance (%)	19.6	13.8	12.4	8.0	7.1	8.4
Cumulative adjusted variance (%)	19.6	33.4	45.8	53.8	60.9	69.3

Example: Pitprops Data

Table 3. Pitprops Data: Loadings of the First Six Sparse PCs by SPCA. Empty cells have zero loadings.

<i>Variable</i>	<i>PC1</i>	<i>PC2</i>	<i>PC3</i>	<i>PC4</i>	<i>PC5</i>	<i>PC6</i>
topdiam	-0.477					
length	-0.476					
moist		0.785				
testsg		0.620				
ovensg	0.177		0.640			
ringtop			0.589			
ringbut	-0.250		0.492			
bowmax	-0.344	-0.021				
bowdist	-0.416					
whorls	-0.400					
clear				-1		
knots		0.013			-1	
diaknot			-0.015			1
Number of nonzero loadings	7	4	4	1	1	1
Variance (%)	28.0	14.4	15.0	7.7	7.7	7.7
Adjusted variance (%)	28.0	14.0	13.3	7.4	6.8	6.2
Cumulative adjusted variance (%)	28.0	42.0	55.3	62.7	69.5	75.8

Example: Pitprops Data

- The SPCA algorithm iterates between the elastic net regression step and the SVD step till convergence.
- The output is the normalized B matrix for $j = 1, \dots, k$

$$\mathbf{v}_j = \hat{\mathbf{b}}_j / \|\hat{\mathbf{b}}_j\|$$

Example: Pitprops Data

- The SPCA algorithm iterates between the elastic net regression step and the SVD step till convergence.
- The output is the normalized B matrix for $j = 1, \dots, k$

$$\mathbf{v}_j = \hat{\mathbf{b}}_j / \|\hat{\mathbf{b}}_j\|$$

Example: Pitprops Data

- The SPCA algorithm iterates between the elastic net regression step and the SVD step till convergence.
- The output is the normalized B matrix for $j = 1, \dots, k$

$$\mathbf{v}_j = \hat{\mathbf{b}}_j / \|\hat{\mathbf{b}}_j\|$$

