Flavors of sparse PCA

Modelization, optimization, and T-Rex 🖜



Arnaud Breloy, 2025



Experiments

Ccl.

Subspace learning grounds

$$\mathbf{x}_i \simeq \mathbf{U}\mathbf{U}^{\top}\mathbf{x}_i$$
, with $\mathbf{U} \in \operatorname{St}(p,k) \stackrel{\Delta}{=} {\mathbf{U} \in \mathbb{R}^{p \times k} \mid \mathbf{U}^{\top}\mathbf{U} = \mathbf{I}}$







Khan data: n = 2308 genes expression of tumors of n = 63 patients, with 4 cancer classes



Khan et al., "Classification and diagnostic prediction of cancers using gene expression profiling and artificial neural networks," Nature medicine, 2001



 $\min_{\mathbf{U} \in \operatorname{St}(p,k)} \quad f(\mathbf{U})$

- **Design** the objective function *f*
- Solve the constrained minimization problem
- Analyze the expected performance, convergence speed, etc.
- Apply the result to some data

Experiments

Ccl.

In this talk

🗸 Design

- Generalize fitting costs from PCA to robust variants
- o Motivating sparse PCA and sparsity promoting penalties

Solve

- $\circ~$ Quick panorama of optimization methods for $\mathrm{St}(p,k)$
- X Analyze
- X Apply

Try something else

 $\circ~$ False discovery rate driven variable selection in SPCA

PCA				
	000		aaaaaaaaaa	

How to design a new "PCA"

Principal component analysis (PCA)

"Vanilla" PCA of rank \boldsymbol{k}

• Singular value decomposition (SVD) of the data matrix $\mathbf{x} = [\mathbf{x}_1, \ \cdots, \ \mathbf{x}_n] \in \mathbb{R}^{p imes n}$

 $\mathbf{X} \stackrel{\mathrm{SVD}}{=} [\mathbf{U} | \mathbf{U}_{\perp}] \, \mathbf{D} \mathbf{V}^{\top}$

- Loading vectors $\mathbf{U} \in \operatorname{St}(p, k)$
- Principal components $\mathbf{x}_i^k = \mathbf{U}^ op \mathbf{x}_i \in \mathbb{R}^k$
- Projected data $\tilde{\mathbf{x}}_i = \mathbf{U}\mathbf{x}_i^k = \mathbf{U}\mathbf{U}^{\top}\mathbf{x}_i$



Experiments

Ccl.

Beyond "vanilla" PCA

Generalizing PCA?

• Find a **model** or **framework** that yields PCA as a solution

 $\min_{\mathbf{U} \in \operatorname{St}(p,k)} \quad f(\mathbf{U})$

• Leverage corresponding **tools** and **extensions**

 $\min_{\mathbf{U} \in \operatorname{St}(p,k)} \quad \mathfrak{f}(\mathbf{U})$

A major motivation: being robust to heavy-tailed distributions and outliers



Experiments

Ccl.

Probabilistic and Bayesian PCA

• Signal plus noise model

$$\mathbf{x} \stackrel{d}{=} \mathbf{U}\mathbf{s} + \mathbf{n}$$

• **Probabilistic PCA** "MLE of **U** defines *f* "

$$\begin{split} \mathbf{s} &\sim \mathcal{N}(\mathbf{0}, \mathbf{\Sigma}_s) \quad \text{and} \quad \mathbf{n} \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I}) \\ &\implies \quad \mathbf{x} \sim \mathcal{N}(\mathbf{0}, \mathbf{U} \mathbf{\Sigma}_s \mathbf{U}^\top + \sigma^2 \mathbf{I}) \end{split}$$

- Extensions by changing assumptions
 - $\circ~$ Elliptical model for $\mathbf{x},\,\mathbf{s},\,\text{or}~\mathbf{n}$
 - $\circ~\ensuremath{\text{Prior distribution on }U}$ directional statistics







• Euclidean distance from \mathbf{x} to $\operatorname{span}(\mathbf{U})$

$$\operatorname{dist}(\mathbf{U}, \mathbf{x}) = \sqrt{\mathbf{x}^{\top} \mathbf{x} - \mathbf{x}^{\top} \mathbf{U} \mathbf{U}^{\top} \mathbf{x}}$$

• Geometric PCA

$$\mathop{\mathrm{minimize}}\limits_{\mathbf{U}\in\mathrm{St}(p,k)}~~\sum_{i=1}^n \mathrm{dist}^2(\mathbf{U},\mathbf{x}_i)$$

• Extensions using alternate distances

• Robust costs:
$$f(\mathbf{U}) = \sum_{i=1}^{n} \rho(\operatorname{dist}^{2}(\mathbf{U},\mathbf{x}_{i}))$$





• Low-rank approximation

 $\begin{array}{ll} \underset{\mathbf{Y}}{\text{minimize}} & ||\mathbf{X} - \mathbf{Y}||_{F}^{2} \\ \text{subject to} & \operatorname{rank}(\mathbf{Y}) = k \end{array}$

- Extensions using alternate decompositions/structures
 - Low-rank plus sparse recovery (Robust PCA)
 - Matrix completion (missing entries)
 - Non-negative matrix factorization
 - Additional structure in the principal components \rightarrow SPCA?

Intro PCA SPCA Optim. Experiments T-Rex Cc aaaa aaaaaa aaaaaa aaaaaa aaaaaaaaaa	ro
---	----

Sparse PCA



Recalling

- Loading vectors $\mathbf{U} = [\mathbf{u}_1, \dots, \mathbf{u}_k] \in \operatorname{St}(p, k)$
- Principal components

$$[\tilde{\mathbf{x}}_i]_j = \mathbf{u}_j^\top \mathbf{z}_i \in \mathbb{R}$$

 $\bullet\,$ The $j^{\rm th}\mbox{-}PC$ is a linear combination of initial variables weighted by the loading vector

Loading vectors should be sparse

- Improved PC's interpretation
- Double duty: dimension reduction & variable selection

Sparse PCA from the optimization point of view

 $\underset{\mathbf{U}\in\mathrm{St}(p,k)}{\text{minimize}} f(\{\mathbf{x}_i\}_{i=1}^n, \mathbf{U}) + \frac{\lambda h(\mathbf{U})}{\lambda h(\mathbf{U})}$

- *f* is a fitting cost cf. many examples
- *h* is a **sparsity-promoting penalty** mimics ℓ_1 or $\ell_{2,1}$ -norm



птата средно страници с

Optimization on St(p, k)

Existing solutions

• Solution #1: Riemannian optimization on St(p, k)

BOU23 Boumal, "An introduction to optimization on smooth manifolds," Cambridge University Press, 2023
ABS09 Absil, Mahony, Sepulchre, "Optimization algorithms on matrix manifolds," Princeton Univ. Press, 2009
EDE98 Edelman, Arias, Smith, "The geometry of algorithms with orthogonality constraints," SIMAX, 1998

• Solution #2: Majorization-Minimization ticks

KIE02 Kiers, "Setting up alternating least squares and iterative majorization algorithms for solving various matrix optimization problems," Computational statistics & data analysis, 2002BRE21 Breloy et al. "MM on the Stiefel Manifold With Application to Robust Sparse PCA", IEEE TSP, 2021

• Solution #3: ADMM splitting tricks

BRE23 Brehier et al. "Robust and globally sparse PCA via MM and variable splitting," ICASSP, 2023UEM19 Uematsu, Fan, Chen, Lv, Lin, "SOFAR: Large-Scale Association Network Learning," IEEE Trans. on IT, 2019

		Optim.
00000	000	

Experiments

T-Rex 00000000000 Ccl.

#1 Riemannian optimization

Tools from differential geometry

- Manifold \mathcal{M} , tangent spaces $T_{\theta}\mathcal{M}$
- Metric $\langle \cdot, \cdot \rangle_{\theta}$
- Retraction $\mathbf{R}_{\theta}: T_{\theta}\mathcal{M} \to \mathcal{M}$

Riemannian gradient descent

 $\langle \operatorname{grad}_{\mathcal{M}} f(\theta), \xi \rangle_{\theta} = \mathrm{D} f(\theta)[\xi]$

 $\theta_{i+1} = R_{\theta_i}(-t_i \operatorname{grad}_{\mathcal{M}} f(\theta_i))$



Pros: elegant, generic and flexible framework, links with information geometryCons: ca be tedious and computationally expensive, no convincing proximal method (yet?)







Pros: super fast (in time), covers a large class of functions

Cons: yet not universal, has to be tailored (or proxy the objective)



xperiments

Ccl.

#2 Majorization-Minimization framework for $\mathrm{St}(p,k)$

Proposal in BRE21

- A systematic trick to deal with orthonormality
- A catalog of surrogates for usual costs
- Applications to probabilistic and sparse PCA

Majorize f on St(p, k) by a linear surrogate

 $g\left(\mathbf{U}|\mathbf{U}^{t}\right) = -2\mathfrak{Re}\left\{\mathrm{Tr}\{\mathbf{R}_{\mathbf{U}^{t}}^{H}\mathbf{U}\}\right\} + \mathrm{const.}$

Minimize g on $St(p, k) \Leftrightarrow$ Projection

 $\mathbf{U}^{t+1} = \mathcal{P}_{ ext{St}} \left\{ \mathbf{R}_{\mathbf{U}^t}^H
ight\}$





#3 (\mathcal{M})-ADMM and distributed algorithms

Variable splitting

$$\begin{array}{ll} \underset{\mathbf{U},\mathbf{V}}{\text{minimize}} & f_{\mathbf{u}}\left(\mathbf{U}\right) + f_{\mathbf{v}}\left(\mathbf{V}\right) \\ \text{subject to} & \mathbf{U} \in \operatorname{St}(p,k) \\ & \mathbf{U} = \mathbf{V} \end{array}$$

Cyclic updates over $\{U,\ V,\ \Gamma\}$ on augmented Lagrangian problem

$$\underset{\{\mathbf{U}_i\}\in\mathrm{St}(p,k),\mathbf{V}}{\text{minimize}} \quad f_{\mathbf{u}}\left(\mathbf{U}\right) + f_{\mathbf{v}}\left(\mathbf{V}\right) + \gamma ||\mathbf{U} - \mathbf{V}||_F^2 + \langle \mathbf{\Gamma}, \mathbf{U} - \mathbf{V} \rangle$$

Pros: generalization to distributed setting, splits problems into simple ones **Cons**: relaxation trick, still requires to handle **U** with previous methods

		Experiments	
aaaaaa			

Some experiments



Experiments

Ccl.

Optimization benchmark

Robust subspace recovery

$$\underset{\mathbf{U}\in \mathrm{St}(p,k)}{\mathrm{minimize}} \quad \sum_{i=1}^{n} \rho_{H}(\mathrm{dist}^{2}(\mathbf{U},\mathbf{x}_{i}))$$

with Huber loss

$$\rho_{\rm H}(t) = \begin{cases} t/\sqrt{T} & \text{if } \leq T \\ 2\sqrt{t} - \sqrt{T} & \text{if } t > T \end{cases}$$

\rightarrow compare iterative algorithms

0000					
Differe	nt algorithms			(and computa	tional bottlenecks
LER17	Quadratic MM, data matr	ix version			rank- k SVD($p \times n$)
	$\mathbf{U}^{t+1} = \mathcal{P}_k\{ ilde{\mathbf{Z}}_t\}, ext{ with } [ilde{\mathbf{Z}}$	$[t]_{:,i} = \sqrt{\rho' (\text{dist})}$	$^{2}\left(\mathbf{U}^{t},\mathbf{z}_{i} ight)\mathbf{\hat{z}}_{i}$		
MAR05	Fixed point heuristic, cov	variance matrix v	ersion		rank- k SVD $(p imes p)$
	$\mathbf{U}^{t+1} = \mathcal{P}_k \left\{ \mathbf{M} \left(\mathbf{U}^t \right) \right\}, \ \mathbf{v}$	with $\mathbf{M}\left(\mathbf{U}^{t} ight)= ilde{\mathbf{Z}}_{t}$	$ ilde{\mathbf{Z}}_t^H$		
MAN02	Steepest descent on Stie	fel			imesthin-SVDs($p imes k$)
	$\mathbf{U}^{t+1} = \mathcal{P}_{\mathrm{St}} \left\{ \mathbf{U}^t + \gamma \nabla_f \right\}$	$\left(\mathbf{U}^{t} ight)ig\},$ with the	right γ		
MAN02	Newton method on Stief	el			$(p \times k)^2$ system
	$\mathbf{U}^{t+1} = \mathcal{P}_{\mathrm{St}} \{ \mathbf{U}^t + \mathbf{Y} \}, \ W$	with $\mathbf{Y} = \operatorname{cpoint}(\mathbf{Y})$	$\mathbf{U}^t, abla_f(\mathbf{U}^t), \mathbf{H}_f(\mathbf{U}^t)$)	
DIN06	Procrustes-MM				thin-SVD($p imes k$)
	$\mathbf{U}^{t+1} = \mathcal{P}_{\mathrm{St}} \left\{ \mathbf{M} \left(\mathbf{U}^{t} \right) \mathbf{U}^{t} \right\}$	}			

.١

Objective value (-optimal value) versus CPU time





Experiments

Average CPU time of an iteration versus size and rank



Experiments

24



Experiments

Ccl.

Robustness and sparsity



• Fitting

- **Least-squares** $\rho(x) = x$
- **Huber Loss** $\rho_H(x)$

- Penalty
 - $\circ \ \ \mathbf{`'PCA''} \ \lambda = 0$
 - **SPCA** $\lambda \neq 0 + h$ is smooth- ℓ_1



Experiments

Ccl.

Subspace recovery

Sampling

$$\{\mathbf{x}_i\}_{i=1}^n = \{\{\mathbf{x}_i^{\text{in}}\}_{i=1}^{n_{\text{in}}}, \{\mathbf{x}_i^{\text{out}}\}_{i=n_{\text{in}}+1}^n\}$$

 $\mathbf{x}^{ ext{in}} ~\sim~ \mathcal{CN}(\mathbf{0}, ext{SNR} imes \mathbf{U}\mathbf{U}^H + \mathbf{I})$

$$\mathbf{x}^{\mathrm{out}} ~~ \sim ~~ \mathcal{CN}(\mathbf{0}, \mathrm{ONR} \times \mathbf{U}_{\perp} \mathbf{U}_{\perp}^{\mathit{H}} + \mathbf{I})$$

Metric

$$\mathrm{AFE} = \mathbb{E}\left[\mathrm{Tr}\{\hat{\mathbf{U}}^H\mathbf{U}\mathbf{U}^H\hat{\mathbf{U}}\}/k\right]$$



p = 100, k = 10, SNR = 10, dense loadings



Experiments

Ccl.

Subspace recovery

Sampling

$$\{\mathbf{x}_i\}_{i=1}^n = \{\{\mathbf{x}_i^{\text{in}}\}_{i=1}^{n_{\text{in}}}, \{\mathbf{x}_i^{\text{out}}\}_{i=n_{\text{in}}+1}^n\}$$

 $\mathbf{x}^{ ext{in}} ~\sim~ \mathcal{CN}(\mathbf{0}, ext{SNR} imes \mathbf{U}\mathbf{U}^H + \mathbf{I})$

$$\mathbf{x}^{\mathrm{out}} ~~ \sim ~~ \mathcal{CN}(\mathbf{0}, \mathrm{ONR} \times \mathbf{U}_{\perp} \mathbf{U}_{\perp}^{\mathit{H}} + \mathbf{I})$$

Metric

$$AFE = \mathbb{E}\left[\mathrm{Tr}\{\hat{\mathbf{U}}^{H}\mathbf{U}\mathbf{U}^{H}\hat{\mathbf{U}}\}/k\right]$$



n, p = 100, k = 10, SNR = 10, sparse loadings

	000

Experiments

Ccl.

Robustness to outliers







		T-Rex	

FDR controlled SPCA

Experiments

T-Rex

Ccl.

Sparse PCA: EV vs SP trade-off

Limitation

- SPCA is only driven by **explained variance** (EV)
- Not always relevant for variable selection
- Might capture high EV outliers

New point of view in this work

- Links with SPCA
- Focuses on variable selection step
- Reformulates the optimization with promising criterion

FDR = "% of allowed false variable selection"





Inspiration:

Zou, Hastie, Tibshirani, "Sparse principal component analysis," J. of computational and graphical statistics, 2006

2-step SPCA given pre-computed plug-in PCs $\tilde{\mathbf{X}}$

$$\underset{\mathbf{U}\in\operatorname{St}(p,k)}{\operatorname{minimize}} \quad ||\tilde{\mathbf{X}}-\mathbf{U}^{\top}\mathbf{X}||_F^2 + \lambda ||\mathbf{U}||_1$$

Relaxation of orthogonality constraint \rightarrow series of k **elastic-net** problems

$$\underset{\boldsymbol{\beta}_j}{\text{minimize}} \quad ||[\tilde{\mathbf{X}}]_{j,:}^\top - \mathbf{X}^\top \boldsymbol{\beta}_j||_F^2 + \lambda ||\boldsymbol{\beta}_j||^2 + \lambda_1 ||\boldsymbol{\beta}_j||_1$$

Sparse loading vectors $\mathbf{U}_{SPCA} = [\boldsymbol{\beta}_1, \ldots, \boldsymbol{\beta}_k]$



SPCA boils down to solve k problems

$$\underset{\boldsymbol{\beta}_j}{\text{ninimize}} \quad ||[\tilde{\mathbf{X}}]_{j,:}^{\top} - \mathbf{X}^{\top} \boldsymbol{\beta}_j||_F^2 + \lambda ||\boldsymbol{\beta}_j||^2 + \lambda_1 ||\boldsymbol{\beta}_j||_1$$

and tune λ_1 to achieve a **good SP-vs-EV**

Can we change the paradigm?

- move away from this trade-off
- have guarantees on false-discovery-rate
- be more computationally efficient

We are going to ask all of this to the $\ensuremath{\text{T-Rex}}$ %



False Discovery Rate (FDR) and True Positive Rate (TPR)

- # samples: n
- # candidate variables: p > n
- Active variables: $\mathcal{A} \subseteq \{1, \dots, p\} \quad \leftarrow$ the "true" support of our loadings
- Selected variables: $\widehat{\mathcal{A}} \subseteq \{1, \dots, p\}$

$$\mathrm{FDR} := \mathbb{E}\big[\mathrm{FDP}\,\big] := \mathbb{E}\Big[\frac{|\widehat{\mathcal{A}} \setminus \mathcal{A}|}{1 \vee |\widehat{\mathcal{A}}|}\Big] \qquad \text{and} \qquad \mathrm{TPR} := \mathbb{E}\big[\mathrm{TPP}\,\big] := \mathbb{E}\Big[\frac{|\mathcal{A} \cap \widehat{\mathcal{A}}|}{1 \vee |\mathcal{A}|}\Big]$$

Goal

 $\label{eq:eq:alpha} \begin{array}{ll} \max_{\widehat{\mathcal{A}}} \mbox{ TPR} & \mbox{ s.t. } \mbox{ FDR} \leq \alpha, \end{array}$

where $\alpha \in [0, 1]$ is the user-defined target FDR level.

Terminating-Random Experiments (T-Rex): the ingredients

- LARSalgorithm
 - Solves LASSO or Elastic-Net iteratively
 - Produces a sequence of selected variables
- Terminated random experiment
 - \circ Append L dummies to $\mathbf{X} := [\mathbf{X}; \mathbf{X}]$ usually $\sim \mathcal{N}(0, 1)$
 - $\circ\,$ Apply LARS and stop when $T\,{\rm dummies}$ are selected

• T-Rex selector

- Run term. rand. exp. in parallel and vote for \hat{A} threshold $v \in (0.5, 1]$
- FDP can be estimated at this step
- $\circ~$ Calibrate $~T~{\rm and}~v~{\rm to}~{\rm maximize}~{\rm TPP}~{\rm and}~{\rm achieve}~{\rm desired}~{\rm FDR}~{\rm target}$

T-Rex

PCA

SPCA

Optim.

Experiments

T-Rex

Ccl.

T-Rex Selector



• $L \in \mathbb{N}_+$ generated dummies • $T \in \{1, \dots, L\}$ included dummies before stopping

• $v \in [0.5, 1)$ voting level • $\Phi_{T,L}(j)$ relative occurrence of variable j • $\widehat{\mathcal{A}}_L(v, T) := \{j : \Phi_{T,L}(j) > v\}$



Zou, Hastie, Tibshirani, "Sparse principal component analysis," J. of computational and graphical statistics, 2006

2-step SPCA given pre-computed plug-in PCs $\tilde{\mathbf{X}}$

 $\underset{\mathbf{U}\in\operatorname{St}(p,k)}{\operatorname{minimize}} \quad ||\tilde{\mathbf{X}}-\mathbf{U}^{\top}\mathbf{X}||_{F}^{2}+\lambda||\mathbf{U}||_{1}$

Relaxation of orthogonality constraint \rightarrow series of k **elastic-net** problems



Sparse loading vectors $\mathbf{U}_{SPCA} = [\boldsymbol{\beta}_1, \ldots, \boldsymbol{\beta}_k]$

		T-Rex

Cleaning PCA with the T-Rex Selector

Setup (see setup in R pack. scPCA): 4 Classes, n = 60, p = 150, $p_1 = 10$.



(a) FDP = 1, TPP = 1, #Var = 150.



Oracle PCA



T-Rex, tFDR:= 0.3 / 2



(g) FDP = 0, TPP = 1, #Var = 10.







SPCA

Optim.

Experiments

T-Rex

Ccl.

More experiments





Experiments

T-Rex

Ccl.

Experiments on the 28 most influential S&P 500 stocks

WMT COST

AVGC XON

CVX

ARRY

AMZN

AAPL

META

MSFT

JPN

CSCO

BRK-B PG

TSLA

JNJ

PEF MRK NVDA

6006

MA

HD



Raw covariance matrix





PCs removed

PCA DDDDDD SPCA

Optim.

Experiments



Covered so far

🗸 Design

- Generalize fitting costs from PCA to robust variants
- Motivating sparse PCA and sparsity promoting penalties

Solve

 $\circ~$ Quick panorama of optimization methods for $\mathrm{St}(p,k)$

X Analyze

X Apply

Try something else

 $\circ\,$ False discovery rate driven variable selection in SPCA

intro PCA SPCA Optim. Experiments T-Rex CC	Intro
aacao acaoaaa acao acaacao acaocaacaoaaa a	0000

Thanks!



Machkour, J., Breloy, A., Muma, M., Palomar, D. P., & Pascal, F. (2024, April). Sparse PCA with false discovery rate controlled variable selection. In ICASSP 2024-2024 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP) (pp. 9716-9720). IEEE.

Brehier, H., Breloy, A., El Korso, M. N., & Kumar, S. (2023, June). Robust and globally sparse PCA via majorization-minimization and variable splitting. In ICASSP 2023-2023 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP) (pp. 1-5). IEEE.

Collas, A., Bouchard, F., Breloy, A., Ginolhac, G., Ren, C., & Ovarlez, J. P. (2021). Probabilistic PCA from heteroscedastic signals: geometric framework and application to clustering. IEEE Transactions on Signal Processing, 69, 6546-6560.

Breloy, A., Kumar, S., Sun, Y., & Palomar, D. P. (2021). Majorization-minimization on the Stiefel manifold with application to robust sparse PCA. IEEE Transactions on Signal Processing, 69, 1507-1520.

T-Rex demo page https://cran.r-project.org/web/packages/TrexSelector/vignettes/TrexSelector_usage_and_simulations.html