

L1-norm Optimization in Subspace Learning Methods

Presented @ 2012 Korean CVPR Workshop

Nojun Kwak

`nojunk@ajou.ac.kr`

<http://image.ajou.ac.kr>

Department of Electrical Engineering
Ajou University, Korea

Nov. 2, 2012



This presentation is mostly based on the two papers:

- Nojun Kwak, “Principal component analysis based on L1 norm maximization”, IEEE TPAMI, vol. 30, no. 9, pp. 1672 – 1680, Sep. 2008.
- Nojun Kwak and Jiyong Oh, “Feature extraction for one-class classification problem: Enhancements to biased discriminant analysis”, Pattern Recognition, vol. 42, no. 1, pp. 17 – 26, Jan. 2009.



- Introduction
 - L2-PCA: Two interpretations
 - Other flavors of PCA (L1-PCA, R1-PCA)
- PCA-L1
 - PCA-L1: Problem formulation
 - PCA-L1: Algorithm
 - PCA-L1: Examples & Experimental Results
- L1-BDA
 - L1-BDA: Application of the Theorem to BDA
 - L1-BDA: Experimental Results
- Conclusions and future works



- PCA (Principal Component Analysis) [1]
 - Dimensionality reduction technique
 - Data visualization
 - Face recognition (eigenface)
 - A lot of applications
- Pros and Cons of PCA
 - Computationally efficient (SVD)
 - Prone to outliers
 - Instead of L2-norm, L1-norm is used here.



- $X = [\mathbf{x}_1, \dots, \mathbf{x}_n] \in \mathbb{R}^{d \times n}$: dataset
 - d : dimension of input space
 - n : number of samples
 - $\{\mathbf{x}_i\}_{i=1}^n$ is assumed to have zero mean.
- $W \in \mathbb{R}^{d \times m}$: projection matrix
 - m : dimension of feature space (no. of features to be extracted)
 - $\{\mathbf{w}_k\}_{k=1}^m$: set of m projection vectors
- $V \in \mathbb{R}^{m \times n}$: coefficient matrix
 - $V = W^T X$
 - v_{ij} : i th coordinate of \mathbf{x}_j after projection
- $E = X - WV = (I_d - WW^T)X$: error matrix in the original input space



Two interpretations of the conventional L2-PCA I

- 1 Minimizing the error of the projection (in L2-norm)

$$W^* = \underset{W}{\operatorname{argmin}} E_2(W) \quad (1)$$

subject to $W^T W = I_m$, where,

$$\begin{aligned} E_2(W) &= \|X - WV\|_F^2 = \sum_{i=1}^n \left\| \mathbf{x}_i - \sum_{k=1}^m \mathbf{w}_k v_{ki} \right\|_2^2 \\ &= \sum_{i=1}^n \sum_{j=1}^d (x_{ji} - \sum_{k=1}^m w_{jk} v_{ki})^2 \end{aligned} \quad (2)$$



Two interpretations of the conventional L2-PCA II

- ② Maximizing the dispersion of the projection (in L2-norm)

$$W^* = \underset{W}{\operatorname{argmax}} D_2(W) \quad (3)$$

subject to $W^T W = I_m$, where

$$\begin{aligned} D_2(W) &= \sum_{i=1}^n \|W^T \mathbf{x}_i\|_2^2 = \sum_{i=1}^n \sum_{k=1}^m (\mathbf{w}_k^T \mathbf{x}_i)^2 \\ &= \|W^T X\|_F^2 = \operatorname{tr}(W^T S_x W) \end{aligned} \quad (4)$$

$S_x = X X^T$: scatter matrix of X .

- **The two are equivalent!! (Dual)**
→ solved by EVD (of S_x) or SVD (of X).



PCA utilizing other norms than L2 I

- L2 norm (Frobenius) is prone to outliers.
- Instead of L2, use other norm in the optimization
- ① Minimization of L1 projection error (input space) [2] – [4]

$$W^* = \underset{W}{\operatorname{argmin}} E_1(W) \quad \text{subject to} \quad WW^T = I_m. \quad (5)$$

$$\begin{aligned} E_1(W) &= \|X - WV\|_1 = \sum_{i=1}^n \left\| \mathbf{x}_i - \sum_{k=1}^m \mathbf{w}_k v_{ki} \right\|_1 \\ &= \sum_{i=1}^n \sum_{j=1}^d |x_{ji} - \sum_{k=1}^m w_{jk} v_{ki}|. \end{aligned} \quad (6)$$

- **Not invariant to rotations** of the input space.
- Exact solution of (5) is hard to achieve. (iterative solutions)



- 2 Minimization of R1-norm of the projection error [5]

$$W^* = \underset{W}{\operatorname{argmin}} E_{R1}(W) \quad \text{subject to} \quad WW^T = I_m. \quad (7)$$

$$E_{R1}(W) = \|X - WV\|_{R1} \triangleq \sum_{i=1}^n \left(\sum_{j=1}^d (x_{ji} - \sum_{k=1}^m w_{jk}v_{ki})^2 \right)^{\frac{1}{2}}. \quad (8)$$

- R1-norm is a hybrid of L1 and L2.
- **Optimal solution depends on m** i.e., $W^{1*} \neq W^{2*}$.
- Solution is **not intuitive, not exact**. (Huber's M-estimator is used.)



Formulation of PCA-L1 [6] I

- Motivation:
 - Previous methods (L1-PCA, R1-PCA) minimizes projection error (E , 1st interpretation).
 - Instead of solving minimization problem, **maximize the dispersion of projection** (D , 2nd interpretation).
- Problem formulation

$$W^* = \underset{W}{\operatorname{argmax}} D_1(W) \quad \text{subject to} \quad W^T W = I_m \quad (9)$$

$$D_1(W) = \sum_{i=1}^n \|W^T \mathbf{x}_i\|_1 = \sum_{i=1}^n \sum_{k=1}^m |w_k^T \mathbf{x}_i| \quad (10)$$

- $W^T W = I_m$: to ensure orthonormality of the projection vectors. (not necessary)
- (9) and (5) are **not equivalent!**



Formulation of PCA-L1 [6] II

- Pros and Cons of (9)
 - (9) are invariant to rotations.
 - As R1-PCA, the solution depends on m .
- Modified problem: $m = 1$

$$\mathbf{w}^* = \underset{\mathbf{w}}{\operatorname{argmax}} \|\mathbf{w}^T X\|_1 = \underset{\mathbf{w}}{\operatorname{argmax}} \sum_{i=1}^n |\mathbf{w}^T \mathbf{x}_i| \quad (11)$$

subject to $\|\mathbf{w}\|_2 = 1$.



Algorithm: PCA-L1

- 1 Initialization: Pick any $\mathbf{w}(0)$. Set $\mathbf{w}(0) \leftarrow \mathbf{w}(0)/\|\mathbf{w}(0)\|_2$ and $t = 0$.
- 2 Polarity check: For all $i \in \{1, \dots, n\}$, if $\mathbf{w}^T(t)\mathbf{x}_i < 0$, $p_i(t) = -1$, otherwise $p_i(t) = 1$.
- 3 Flipping and maximization: Set $t \leftarrow t + 1$ and $\mathbf{w}(t) = \sum_{i=1}^n p_i(t-1)\mathbf{x}_i$. Set $\mathbf{w}(t) \leftarrow \mathbf{w}(t)/\|\mathbf{w}(t)\|_2$.
- 4 Convergence check:
 - a. If $\mathbf{w}(t) \neq \mathbf{w}(t-1)$, go to Step 2.
 - b. Else if there exists i such that $\mathbf{w}^T(t)\mathbf{x}_i = 0$, set $\mathbf{w}(t) \leftarrow (\mathbf{w}(t) + \Delta\mathbf{w})/\|\mathbf{w}(t) + \Delta\mathbf{w}\|_2$ and go to Step 2. Here, $\Delta\mathbf{w}$ is a small nonzero random vector.
 - c. Otherwise, set $\mathbf{w}^* = \mathbf{w}(t)$ and stop.



Optimality of the Algorithm I

Theorem

With the above PCA-L1 procedure, the projection vector \mathbf{w} converges to \mathbf{w}^* , which is a local maximum point of $\sum_{i=1}^n |\mathbf{w}^T \mathbf{x}_i|$.

Proof.

Firstly, we can show that $\sum_{i=1}^n |\mathbf{w}^T(t) \mathbf{x}_i|$ is a non-decreasing function of t as the following:

$$\begin{aligned} \sum_{i=1}^n |\mathbf{w}^T(t) \mathbf{x}_i| &= \mathbf{w}^T(t) \left(\sum_{i=1}^n p_i(t) \mathbf{x}_i \right) \geq \mathbf{w}^T(t) \left(\sum_{i=1}^n p_i(t-1) \mathbf{x}_i \right) \\ &\geq \mathbf{w}^T(t-1) \left(\sum_{i=1}^n p_i(t-1) \mathbf{x}_i \right) = \sum_{i=1}^n |\mathbf{w}^T(t-1) \mathbf{x}_i|. \end{aligned} \tag{12}$$

$\therefore \mathbf{w}(t) \rightarrow \mathbf{w}^*$ in finite number of iterations. □



Proof.

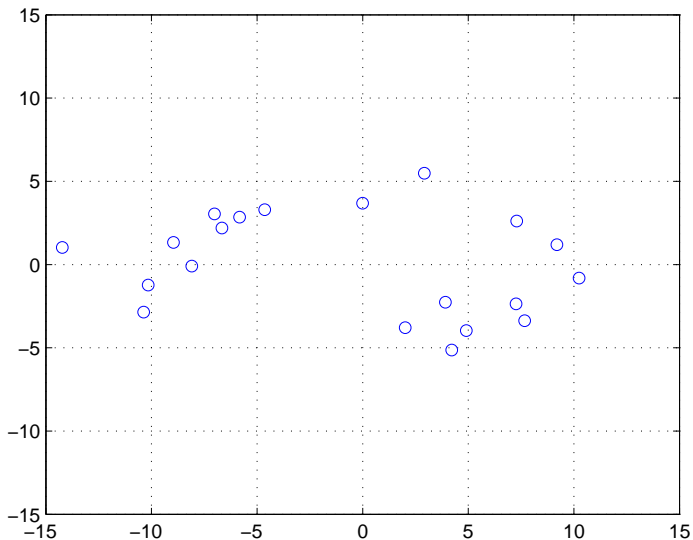
Local maximality of \mathbf{w}^* :

- 1 $\mathbf{w}^{*T} p_i(t) \mathbf{x}_i \geq 0 \forall i. \because \mathbf{w}(t)$ converges to \mathbf{w}^*
- 2 \exists a small neighbor $N(\mathbf{w}^*)$ of $\mathbf{w}^* \ni$ if $\mathbf{w} \in N(\mathbf{w}^*)$, then $\mathbf{w}^T p_i(t) \mathbf{x}_i \geq 0 \forall i$.
 - Finite number of data points.
 - $\mathbf{w}^{*T} \mathbf{x}_i \neq 0 \forall i$. (Step 4b)
- 3 $\sum_{i=1}^n |\mathbf{w}^{*T} \mathbf{x}_i| > \sum_{i=1}^n |\mathbf{w}^T \mathbf{x}_i| \forall \mathbf{w} \in N(\mathbf{w}^*)$
 - $\because \mathbf{w}^* \parallel \sum_{i=1}^n p_i(t) \mathbf{x}_i$
 - $\rightarrow \mathbf{w}^*$ is a local maximum point.

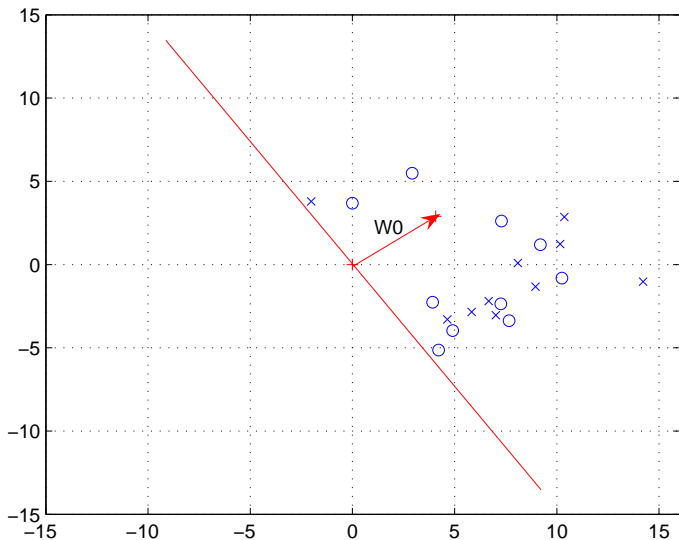
Therefore, the PCA-L1 procedure finds a local maximum point \mathbf{w}^* . □



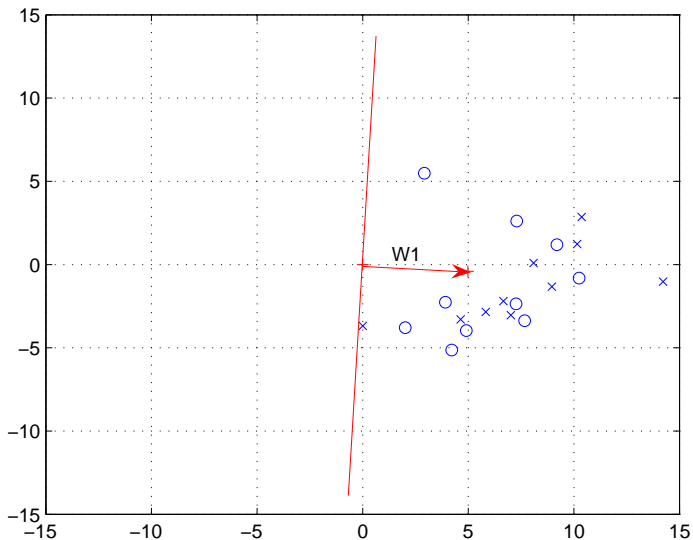
Simple Example (2D case)



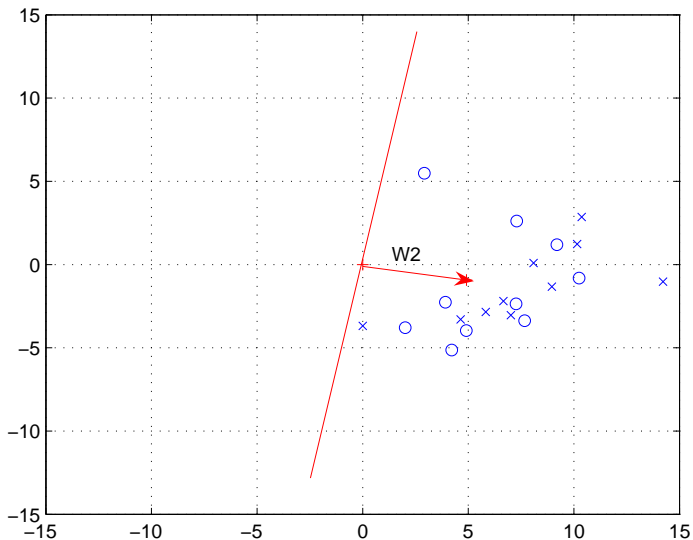
Simple Example (2D case)



Simple Example (2D case)



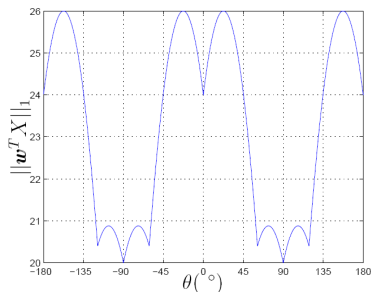
Simple Example (2D case)



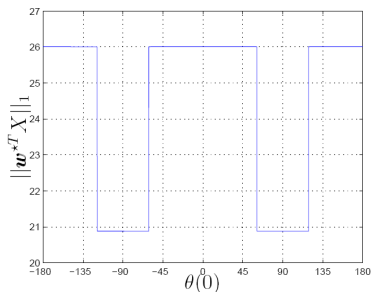
Dependency on Initial Projection Vector $w(0)$

- An example with 5 data points (in 2D)

$$X = \begin{bmatrix} 0 & 9 & -9 & 3 & -3 \\ 10 & -5 & -5 & 0 & 0 \end{bmatrix}.$$



(a) Objective function



(b) Dependency on initial vector $w(0)$ ($\theta(0)$)

- Final vector w^* depends on initial point $w(0)$.
 - Start with various initial points.
 - Set $w(0) = w_{L2-PCA}$.



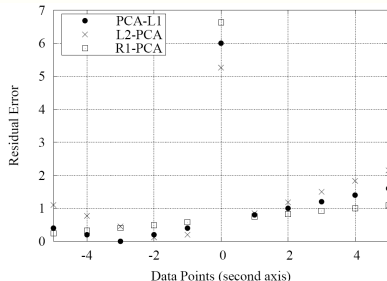
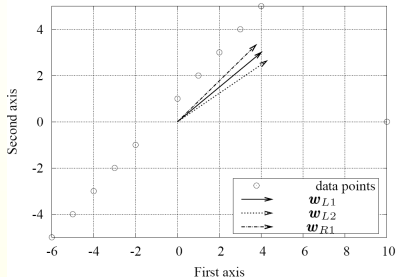
(Greedy search algorithm)

- 1 $\mathbf{w}_0 = \mathbf{0}, \{\mathbf{x}_i^0 = \mathbf{x}_i\}_{i=1}^n$.
 - 2 For $j = 1$ to m ,
 - a. $\forall i \in \{1, \dots, n\}, \mathbf{x}_i^j = \mathbf{x}_i^{j-1} - \mathbf{w}_{j-1}(\mathbf{w}_{j-1}^T \mathbf{x}_i^{j-1})$.
 - b. In order to find \mathbf{w}_j , apply the PCA-L1 procedure to $X^j = [\mathbf{x}_1^j, \dots, \mathbf{x}_n^j]$.
- Orthogonality is **not necessary**.
 - **However**, to limit the search space, orthogonality is assumed.



Comparison with other versions of PCA

- A toy problem with an outlier



- At a glance, R1-PCA seems to remove the effect of outliers most effectively.
- However, the average residual error is

	L2-PCA	R1-PCA	PCA-L1
Average Residual Error	1.401	1.206	1.200



Classification Results (UCI datasets)

Table: Average Classification Rates on UCI datasets (%). The last column is the averages of the best classification rates among the cases where the number of extracted features was one to half the number of original inputs

No. of extracted features	1	2	3	4	Best performance
L2-PCA	62.49	68.59	73.36	76.79	76.46
R1-PCA	62.44	68.49	73.63	76.52	76.47
PCA-L1	63.94	71.88	74.90	77.43	78.15



Face Reconstruction Results



Figure: Face images with occlusion and the reconstructed faces: 1st column: original, 2nd column: PCA-L1, 3rd column: L2-PCA, 4th column: R1-PCA. (reconstructed with 20 projection vectors)



- BDA (Biased Discriminant Analysis)

- **One-class** classification problem
- Maximizes the ratio between the **positive scatter** and **negative scatter** matrices.

$$W^* = \operatorname{argmax}_W \frac{\operatorname{tr}(W^T S_y W)}{\operatorname{tr}(W^T S_x W)}. \quad (13)$$

$$S_x = \sum_{i=1}^{n_x} (\mathbf{x}_i - \mathbf{m}_x)(\mathbf{x}_i - \mathbf{m}_x)^T \quad (14)$$

$$S_y = \sum_{i=1}^{n_y} (\mathbf{y}_i - \mathbf{m}_x)(\mathbf{y}_i - \mathbf{m}_x)^T, \quad (15)$$

- $\{\mathbf{x}_i\}_{i=1}^{n_x}$: positive samples
- $\{\mathbf{y}_i\}_{i=1}^{n_y}$: negative samples
- \mathbf{m}_x : mean of the positive samples



1 Sphering

a. Solve the EVD problem $S_x U = U \Lambda_1$.

b. Scale each column \mathbf{u}_n of U to make $\{\hat{x}_{in} = \hat{\mathbf{u}}_n^T (\mathbf{x}_i - \mathbf{m}_x)\}_{i=1}^{n_x}$ have unit variance.

$\hat{\mathbf{u}}_n$: a scaled version of \mathbf{u}_n

If $\|\mathbf{u}_n\| = 1$, this is equivalent to setting $\hat{\mathbf{u}}_n = \mathbf{u}_n / \sqrt{\lambda_{1n}}$, where λ_{1n} denotes the n -th eigenvalue of S_x .

2 Maximization

a. Find M weight vectors $\{\mathbf{v}_n\}_{n=1}^M$ that maximizes the following objective function:

$$V = \underset{W}{\operatorname{argmax}} |W^T \hat{S}_y W|, \quad \text{subject to } W^T W = I, \quad (16)$$

where $\hat{S}_y = \hat{U}^T S_y \hat{U}$.

b. Output $W = \hat{U} V$.

Here, $\hat{U} = [\hat{\mathbf{u}}_1, \hat{\mathbf{u}}_2, \dots]$.



Sphering Operation in BDA

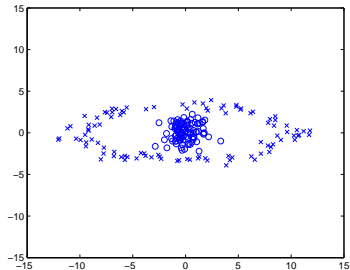
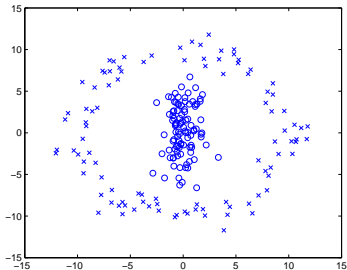


Figure: Data distribution before and after sphering process



- Problem reformulation: L1-BDA

$$W^* = \operatorname{argmax}_W \frac{(\sum_{i=1}^{n_y} |W^T \mathbf{y}_i|)^2}{\operatorname{tr}(W^T S_x W)}. \quad (17)$$

- Solution: 2-step (Sphering \rightarrow Maximization)
 - Replace L2-norm with L1-norm in the Maximization part (after sphering)
 - The L1-norm maximization technique developed for L1-PCA can be directly utilized for maximizing the numerator.



Experimental Results of L1-BDA (FERET Eye Data) I

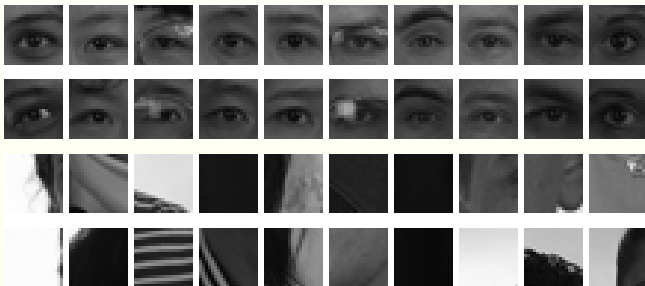


Figure: Eye and noneye samples for training



Experimental Results of L1-BDA (FERET Eye Data) II

Table: Classification rates for FERET Eye dataset. (1-NN)

# of features	LDA	Cher. LDA	BDA	SBDA ($\gamma = 1$)	L1BDA	SL1BDA ($\gamma = 1$)
1	86.875	84.625	81.875	86.25	87.625	86.625
2	–	83.75	88	95	91.125	94.5
3	–	81.375	91.25	96.625	96.375	96.75
4	–	78.75	93	97.5	97.125	97.625
5	–	80.625	93	98.25	97.375	98.25
6	–	79.75	94.875	98.625	97.375	98.625
7	–	78	94.75	98.5	97.875	98.125
8	–	79.125	95	98.125	98	98.125
9	–	78.625	95	98.375	98	98.5
10	–	79.75	95.875	98	97.875	98.125



Experimental Result of L1-BDA (UCI)

Table: Experimental Results on UCI datasets (1-NN classifier)

Data set	LDA	Chernoff LDA	BDA	SBDA ($\gamma = 1$)	L1BDA	SL1BDA ($\gamma = 1$)
Australian	80.55 \pm 1.02 (1)	81.35 \pm 0.82 (6)	80.68 \pm 1.28 (3)	81.29 \pm 0.74 (4)	81.57 \pm 1.19 (4)	81.67\pm1.19 (3)
Balance	88.00 \pm 0.49 (2)	87.68 \pm 0.79 (2)	91.25 \pm 0.99 (3)	95.73 \pm 0.62 (2)	95.76\pm0.64 (2)	94.38 \pm 0.87 (3)
Breast	96.03 \pm 0.34 (1)	96.34\pm0.35 (2)	95.94 \pm 0.44 (2)	96.28 \pm 0.29 (6)	96.05 \pm 0.62 (2)	96.09 \pm 0.20 (6)
Glass	62.85 \pm 1.73 (5)	63.55 \pm 1.43 (9)	66.19 \pm 2.21 (3)	70.65\pm1.12 (6)	67.66 \pm 2.32 (3)	68.74 \pm 1.37 (8)
Heart	76.43 \pm 1.19 (1)	76.60 \pm 1.78 (1)	76.73 \pm 1.41 (2)	76.73 \pm 1.41 (2)	77.44\pm1.41 (4)	77.44\pm1.41 (4)
Iris	96.93 \pm 0.72 (1)	97.13 \pm 0.77 (1)	97.20\pm0.69 (1)	96.87 \pm 0.63 (1)	97.13 \pm 0.77 (1)	96.60 \pm 0.21 (4)
Liver	61.01 \pm 3.28 (1)	65.65\pm2.50 (4)	65.04 \pm 2.20 (3)	65.42 \pm 2.24 (3)	65.13 \pm 1.39 (5)	65.13 \pm 1.39 (5)
Pima	69.13 \pm 1.32 (1)	69.78 \pm 0.51 (8)	69.79 \pm 0.74 (5)	70.01 \pm 0.56 (5)	70.42\pm1.55 (4)	70.42\pm1.55 (4)
Sonar	73.03 \pm 2.27 (1)	81.06 \pm 2.25 (54)	78.56 \pm 2.03 (5)	84.86 \pm 1.59 (18)	80.86 \pm 1.28 (5)	84.90\pm1.51 (17)
Vehicle	74.33 \pm 1.10 (3)	81.55\pm0.51 (9)	73.45 \pm 0.66 (5)	76.32 \pm 0.75 (15)	74.65 \pm 0.46 (2)	80.01 \pm 0.50 (3)
Average	77.83	80.07	79.48	81.41	80.67	81.54



- **PCA-L1**

- finds projections that **maximizes L1-norm in the projected space**.
- is proven to find a **local maximum** point.
- is **robust to outliers**.
- is **simple and easy to implement**.
- is **relatively fast** with small number of iterations.
- The number of iterations **does not depend on the dimension of input space**.
- The same technique can be applied to other feature extraction algorithm.



- Non-greedy version extension

- F. Nie, H. Huang, C. Ding, D. Luo, and H. Wang, "Robust principal component analysis with non-greedy L1-norm maximization", in Proc. 22nd International Conf. on Artificial Intelligence, 2011, pp. 1433 – 1438.

- 2D & tensor extension

- X. Li, Y. Pang, and Y. Yuan, "L1-norm based 2DPCA", IEEE Trans. on SMC-B, vol. 38, no. 4, pp. 1170 – 1175, Aug. 2010.
- Y. Pang, X. Li, and Y. Yuan, "Robust tensor analysis with L1-norm", IEEE Trans. Circuits Syst. Video Techn., vol. 20, no. 2, pp. 172.178, 2010.

- Mean or Median? - Generalized mean

- Jiyong Oh, Nojun Kwak, Minsik Lee and Chong-Ho Choi, "Generalized mean for feature extraction in one-class classification problems", submitted to Pattern Recognition.

- Generalization to Lp-norm

- Nojun Kwak, "Principal component analysis by Lp-norm maximization", submitted to IEEE Trans. on SMC-B.

- Applications to other subspace methods (e.g., LDA)

- Jae Hyun Oh and Nojun Kwak, "Generalization of linear discriminant analysis using Lp-norm", submitted to Pattern Recognition Letters.



- [1] I.T. Jolliffe,
Principal Component Analysis,
Springer-Verlag, 1986.
- [2] A. Baccini, P. Besse, and A.D. Falguerolles,
"A L1-norm pca and a heuristic approach,"
in *Ordinal and Symbolic Data Analysis*, E. Diday, Y. Lechevalier, and P. Opitz, Eds. 1996, pp. 359–368,
Springer.
- [3] Q. Ke and T. Kanade,
"Robust subspace computation using l1 norm,"
Tech. Rep. CMU-CS-03-172, Carnegie Mellon University, Aug. 2003,
<http://citeseer.ist.psu.edu/ke03robust.html>.
- [4] Q. Ke and T. Kanade,
"Robust l1 norm factorization in the presence of outliers and missing data by alternative convex programming,"
in *Proc. IEEE Conference on Computer Vision and Pattern Recognition*, June 2005.
- [5] C. Ding, D. Zhou, X. He, and H. Zha,
"R1-pca: rotational invariant l1-norm principal component analysis for robust subspace factorization,"
in *Proc. International Conference on Machine Learning*, Pittsburgh, PA, June 2006.
- [6] N. Kwak,
"Principal component analysis based on L1 norm maximization,"
IEEE TPAMI, vol. 30, no. 9, pp. 1672–1680, Sep. 2008.
- [7] N. Kwak and J. Oh,
"Feature extraction for one-class classification problem: Enhancements to biased discriminant analysis,"
Pattern Recognition, vol. 42, no. 1, pp. 17–26, Jan. 2009.

