

A Graph Based Approach for Formalizing Subjective Interestingness of Data Projections

Bo Kang¹, Jefrey Lijffijt², Raúl Santos-Rodríguez², Tijl De Bie^{1,2}

¹Ghent University, ²University of Bristol



1. Motivation

- Principal Component Analysis (PCA) is effective: variance in data is often dominated by relevant structure.
- However, not all structure is interesting to every user.
- Subjectively Interesting Component Analysis (SICA) aims to find subjectively interesting projection of the data.

2. Method

- **Dataset:** n data points in d -dimensional real space $\hat{\mathbf{x}} \in \mathbb{R}^d$, represented by means of matrix $\hat{\mathbf{X}} \in \mathbb{R}^{n \times d}$.
- **Prior expectation:** $p_{\mathbf{x}} : \mathbb{R}^{n \times d} \rightarrow \mathbb{R}$ represents prior expectations of user
 - Scale of data:
$$\mathbb{E}_{p_{\mathbf{x}}} \left[\frac{1}{n} \sum_i^n \|\mathbf{x}_i\|^2 \right] = b.$$
- Similarities between pairs from E :

$$\mathbb{E}_{p_{\mathbf{x}}} \left[\frac{1}{|E|} \sum_{(i,j) \in E} \|\mathbf{x}_i - \mathbf{x}_j\|^2 \right] = c.$$

Set E contains all pairs of points expected to be similar, which defines a graph $G([1 \dots n], E)$ over all data points.
- $p_{\mathbf{x}}$ is inferred as MaxEnt distribution subject to these constraints^{1,2}.
- **Subjectively interesting projection:** maximizes information content of event $\hat{\mathbf{X}}\mathbf{W} \in [\hat{\Pi}_{\mathbf{W}}, \hat{\Pi}_{\mathbf{W}} + \Delta \mathbf{I}]$ with projection $\hat{\Pi}_{\mathbf{W}} \triangleq \hat{\mathbf{X}}\mathbf{W} \in \mathbb{R}^{n \times k}$.

$$\max_{\mathbf{W}} \text{Tr}(\mathbf{W}'\hat{\mathbf{X}}'[\lambda\mathbf{I} + \mu\mathbf{L}]\hat{\mathbf{X}}\mathbf{W}),$$

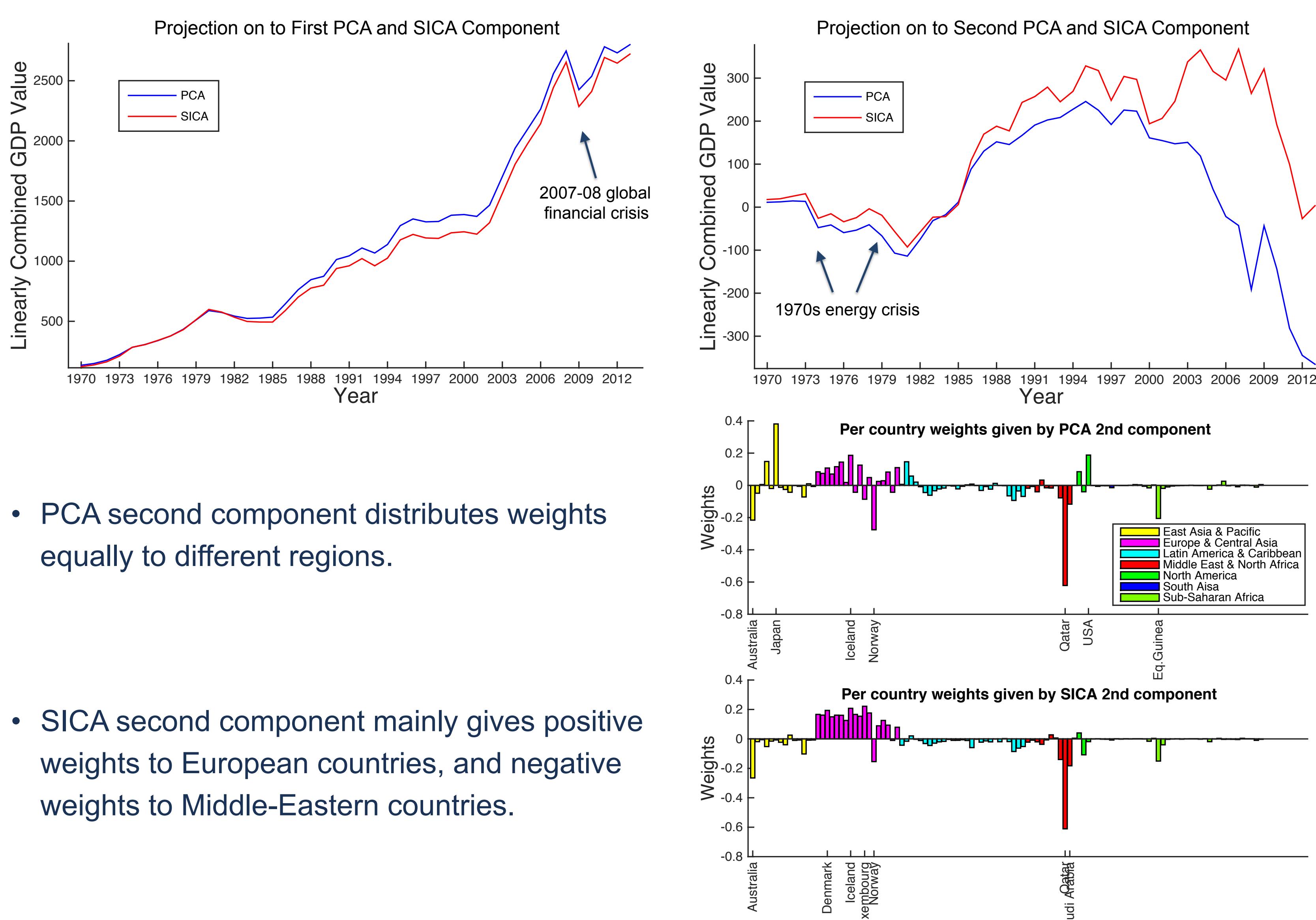
s.t. $\mathbf{W}'\mathbf{W} = \mathbf{I}$.

¹Tijl De Bie. An information theoretic framework for data mining. In Proc. of KDD, pages 564–572, 2011.

²Tijl De Bie. Subjective interestingness in exploratory data mining. In Proc. of IDA, pages 19–31, 2013.

3.2 Case Study: Time Series

- **Dataset:** GDP per Capita⁵ of 110 countries between year 1970 and 2013, i.e., $\hat{\mathbf{X}} \in \mathbb{R}^{44 \times 110}$. Countries are categorized into seven regions.
- **Prior expectation:** GDP value between adjacent years is unlikely to have drastic change. The resulted graph is a chain with 44 nodes.
- **Results:** projections on to first PCA and SICA components show a rather smooth increase over years.
- The projection on to second SICA component shows more local fluctuations.

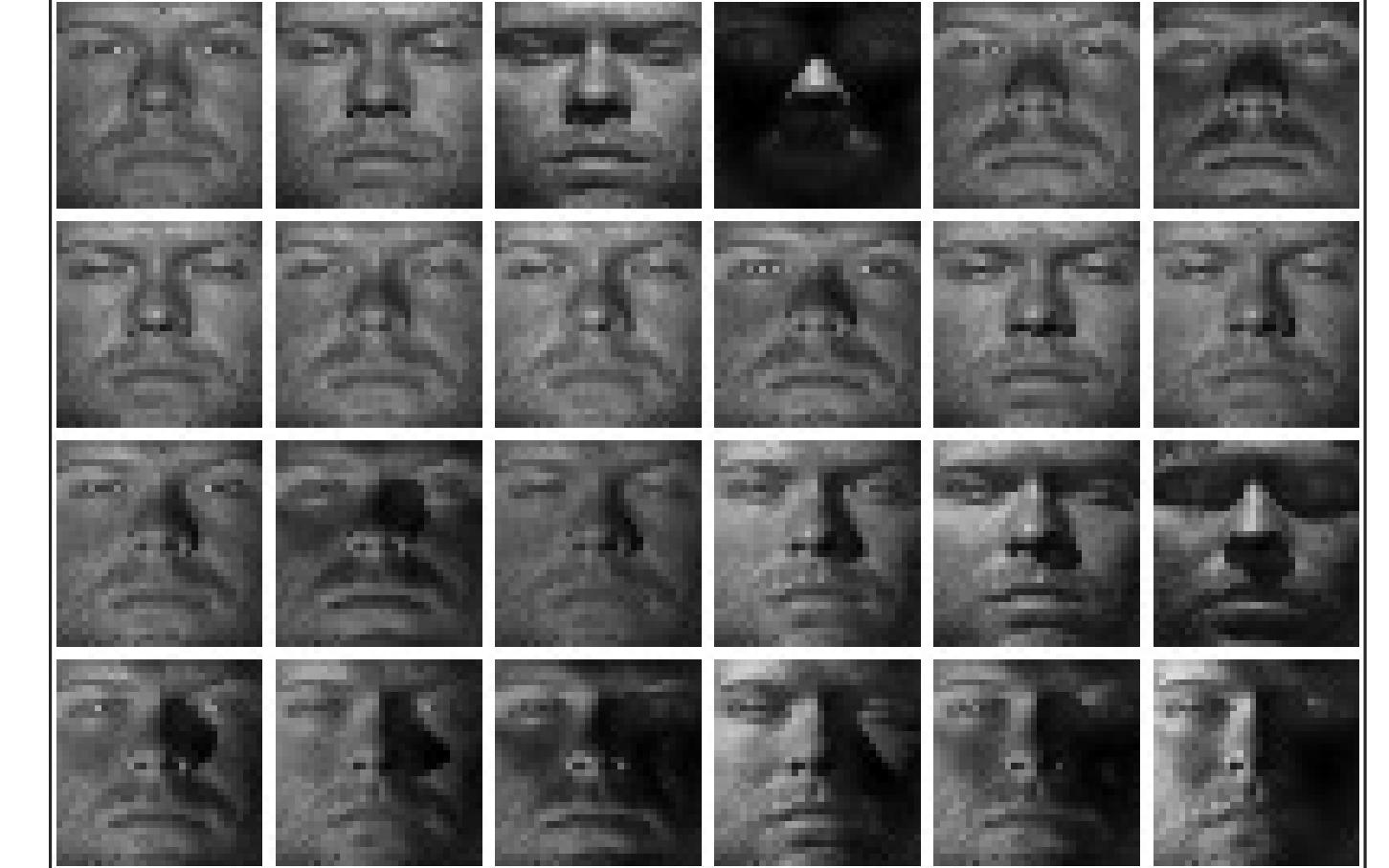


- PCA second component distributes weights equally to different regions.

- SICA second component mainly gives positive weights to European countries, and negative weights to Middle-Eastern countries.

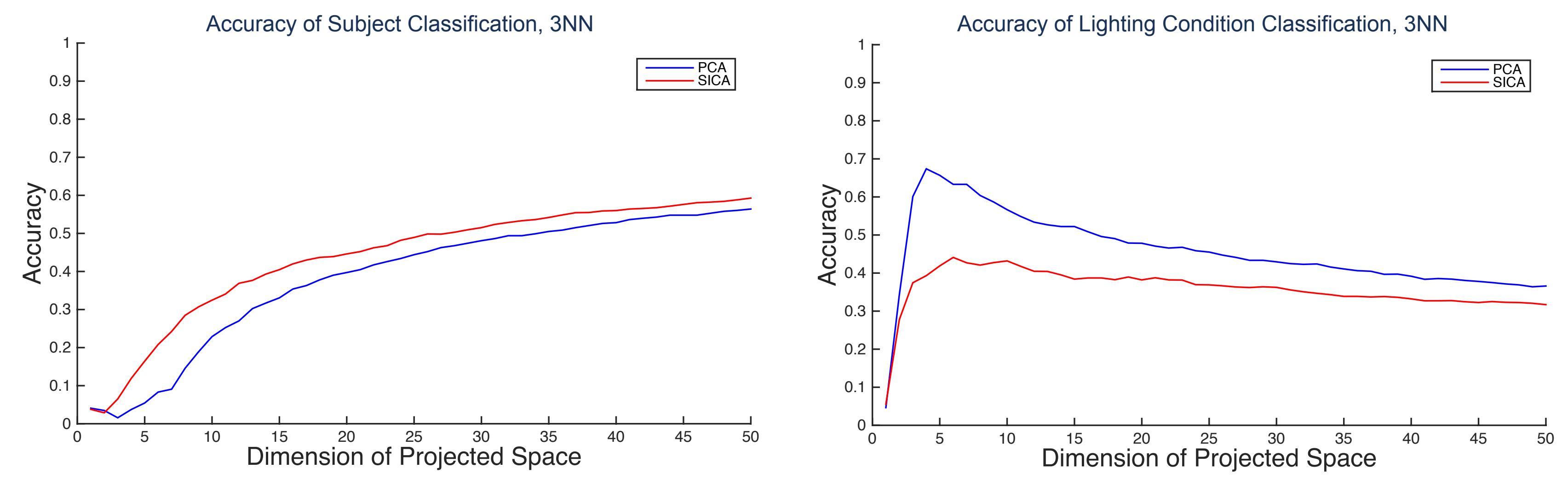
3.1 Case Study: Images

- **Dataset:** 1684 gray-scale frontal images (32x32 pixels) of 31 human subjects under 64 lighting conditions, compiled from the Extended Yale Face Database B^{3,4}, i.e., $\hat{\mathbf{X}} \in \mathbb{R}^{1684 \times 1024}$.

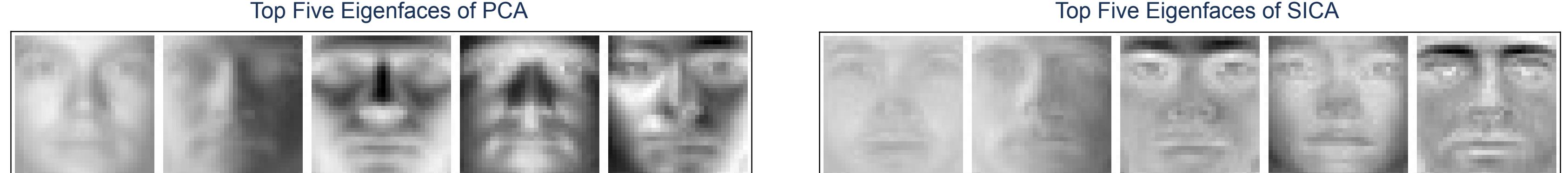


Subject One, First 24 Lighting Conditions

- **Prior expectation:** images with same lighting are similar. Resulting graph constraint consists of 64 cliques.
- **Results:** apply 3-NN on subject classification in projected feature space, SICA achieved better classification accuracy, while on lighting condition classification, PCA achieved better accuracy.



- The top PCA eigenfaces reflect lighting conditions, while the top SICA ones reflect facial features.

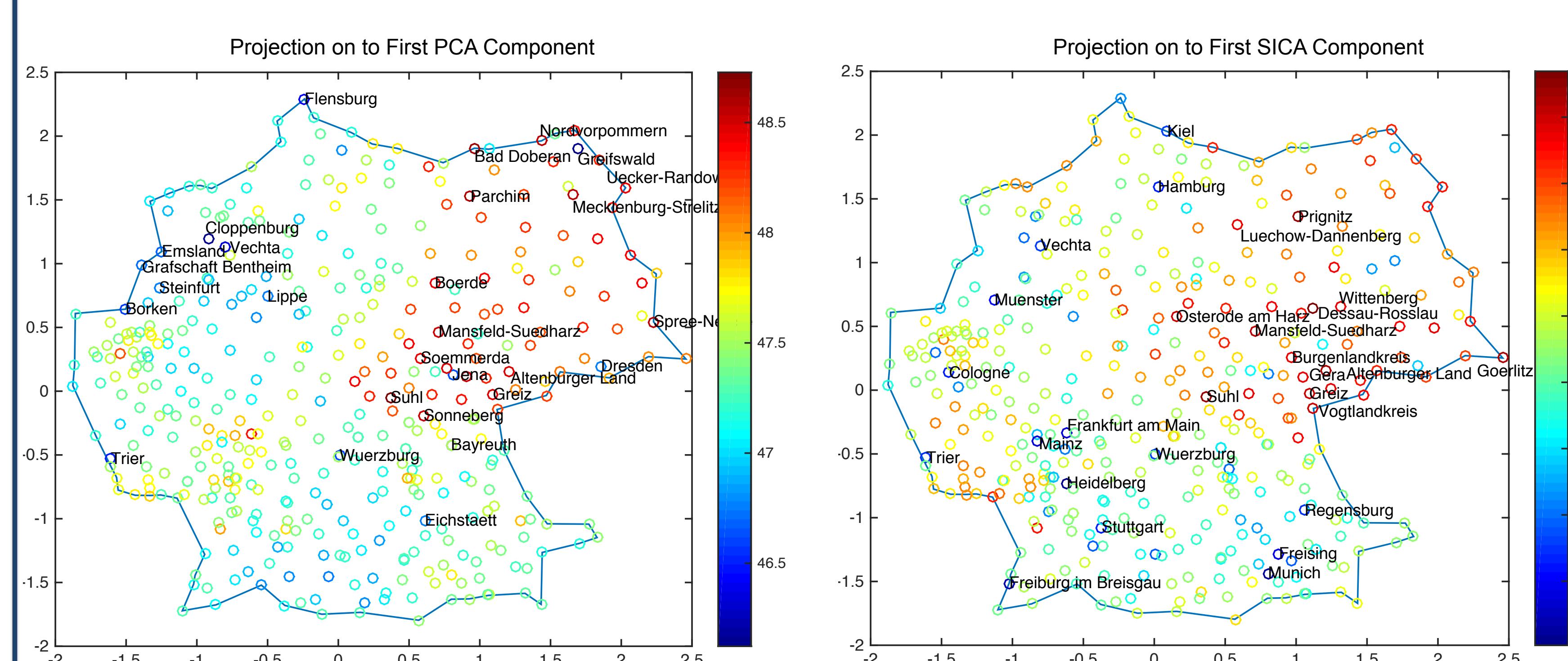


³A.S. Georghiades, P.N. Belhumeur, and D.J. Kriegman. From few to many: illumination cone models for face recognition under variable lighting and pose. *IEEE Trans. Pattern Anal. Mach. Intelligence*, 23 (6): 643–660, 2001.

⁴Kuang-Chih Lee, Jeffrey Ho, and David J Kriegman. Acquiring linear subspaces for face recognition under variable lighting. *Pattern Analysis and Machine Intelligence, IEEE Transactions on*, 27(5):684–698, 2005.

3.3 Case Study: Spatial Data

- **Dataset:** age structure of 412 districts (Landkreise) in Germany⁶. It contains five categories: Elder (age > 64), Old (between 45 and 64), Middle Aged (between 25 and 44), Young (between 18 and 24), Children (age < 18), i.e., $\hat{\mathbf{X}} \in \mathbb{R}^{412 \times 5}$.
- **Prior expectation:** historically, population density and birth rate in eastern Germany are lower than the rest of the country. This information corresponds to a graph constraint with two cliques.
- **Results:** data points are colored by their weighted scores w.r.t the first PCA and SICA component.



- Interpret results by inspecting the elements of the first PCA and SICA component.

	Elder	Old	Mid-Age	Young	Children
PCA First Component	0.44	0.60	0.54	0.17	0.35
SICA First Component	0.62	0.31	-0.69	-0.19	-0.05