

Criteria for choosing the number of dimensions in a principal component analysis: An empirical assessment

Renata B. Silva¹, Daniel de Oliveira², Davi P. Santos³,
Lucio F. D. Santos⁴, Rodrigo E. Wilson¹, and Marcos Bedo¹

¹Fluminense Federal University (INFES/UFF), Brazil

{renatabarbosa, rodrigoerthal, marcosbedo}@id.uff.br

²Fluminense Federal University (IC/UFF), Brazil

danielcmo@ic.uff.br

³University of São Paulo (ICMC/USP), Brazil

davips@icmc.usp.br

⁴Federal Institute of North of Minas Gerais (IFNMG), Brazil

lucio.santos@ifnmg.edu.br

Abstract. *Principal component analysis (PCA) is an efficient model for the optimization problem of finding d' axes of a subspace $\mathbb{R}^{d'} \subseteq \mathbb{R}^d$ so that the mean squared distances from a given set \mathcal{R} of points to the axes are minimal. Despite being steadily employed since 1901 in different scenarios, e.g., mechanics, PCA has become an important link in machine learning chained tasks, such as feature learning and AutoML designs. A frequent yet open issue that arises from supervised-based problems is how many PCA axes are required for the performance of machine learning constructs to be tuned. Accordingly, we investigate the behavior of six independent and uncoupled criteria for estimating the number of PCA axes, namely Scree-Plot %, Scree Plot Gap, Kaiser-Guttman, Broken-Stick, ρ -Score, and 2D. In total, we evaluate the performance of those approaches in 20 high dimensional datasets by using (i) four different classifiers, and (ii) a hypothesis test upon the reported F-Measures. Results indicate Broken-Stick and Scree-Plot % criteria consistently outperformed the competitors regarding supervised-based tasks, whereas estimators Kaiser-Guttman and Scree-Plot Gap delivered poor performances in the same scenarios.*

1. Introduction

Principal component analysis (PCA) is a widely adopted model for dimensionality reduction¹, a pre-processing step related to machine learning tasks [Pearson 1901, Aggarwal 2015]. Such a step is particularly relevant for supervised-driven problems, in which the *curse of dimensionality* [Pestov 2008] may disrupt the learning bias of certain classifiers, e.g., Naïve-Bayes (NB), Instance-based Learning (IbL), Decision-Tree (DT), and Multi-Layer Perceptron (MLP), as well severely degraded their computational performance [Aggarwal 2015, James et al. 2013]. Formally, given a dataset $\mathcal{R} \subset \mathbb{R}^d$, PCA

¹The most relevant *dimensions* for a particular set of points are the most prominent data *features*. Accordingly, we use the terms dimensions and features alternately.

enables finding the d' orthogonal axes of a subspace $\mathbb{R}^{d'} \subseteq \mathbb{R}^d$ so that the mean squared distances from elements in \mathcal{R} to the axes are minimal.

A common yet open issue that arises in practice is distinguishing *relevant* and *non-relevant* axes so that data are reduced to a proper subspace [Pestov 2008, Aggarwal 2015]. Unlike previous approaches that investigate the relationship between d' and co-variance patterns within artificial data [Jackson 1993, Neto et al. 2005], we focus on examining distinct criteria for choosing the number of PCA axes whose performance is assessed by different classifiers. Figure 1 highlights the challenges of exhaustively estimating number d' in supervised UCI dataset WINE² regarding two wrapped classifiers: (i) labeling performance is not monotonic with d' , and (ii) individual maxima are overfitting-prone.

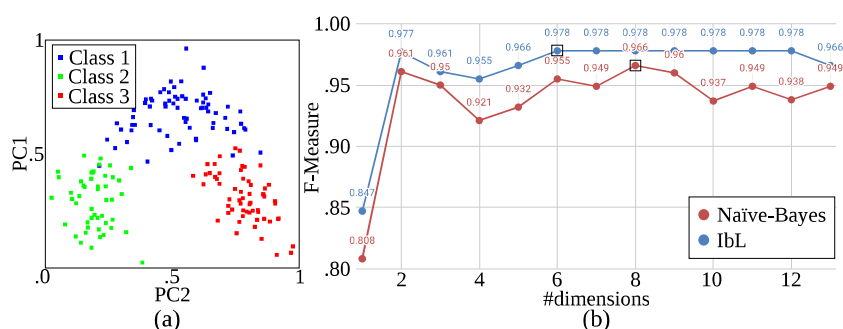


Figure 1. PCA reduction and labeling of WINE dataset.

In this study, we investigate the behavior of six global, distinct, and classifier-unwrapped criteria for choosing the number of dimensions in PCA reductions for labeling problems, namely (i) graphical-based estimators Scree-Plot %, Scree Plot Gap, and 2D; (ii) statistical-based indicators Broken-Stick, and Kaiser-Guttman; and (iii) intrinsic dimension-based criterion ρ -Score. We compared those criteria in the labeling of 20 datasets with four classifiers (NB, IbL, DT, and MLP), and results indicate estimators Broken-Stick and Scree-Plot % surpassed the competitors, while indicators Kaiser-Guttman and Scree-Plot Gap performed modestly. Such outcomes provide indications to devise the tuning of PCA-based pieces within AutoML designs.

This remainder of this paper is organized as follows. Section 2 discusses the estimators for the number of PCA dimensions, while Section 3 describes the material and methods. Sections 4 and 5 provide the experimental comparison and conclude the study.

2. Preliminaries

PCA axes can be found as the uncorrelated coefficients calculated from discrete data features. In a nutshell, reducing a dataset $\mathcal{R} \subset \mathbb{R}^d$ by PCA into a d' -dimensional representation, $d' \leq d$, is a sequence of six sequential steps, namely: (i) scale each \mathcal{R} feature to the $[0, 1]$ interval, (ii) calculate means $\mu_{i, i \in [1, d]}$ for every \mathcal{R} feature, (iii) subtract means $\mu_{i, i \in [1, d]}$ from each \mathcal{R} element, (iv) calculate co-variance matrix $\mathcal{C}_{d \times d}$ from \mathcal{R} entries, (v) obtain both \mathcal{C} eigenvalues and eigenvectors, and (vi) calculate the cross product of \mathcal{R} entries and d' eigenvectors related to the d' *highest and descending-sorted* eigenvalues.

²Data links at github.com/Renata-Barbosa/cpca

Existing criteria for estimating the d' value can be divided into (i) graphical, (ii) statistical, and (iii) intrinsic dimension-based approaches [Neto et al. 2005, Pestov 2008]. Representative approaches of the first group include the following rules:

Scree-Plot % (SP-%). The estimator assumes a few eigenvalues concentrate the largest part of data variance, and calculates d' so that a *percentage* of variance is kept. Figure 2(a) illustrates the *SP-%* rule for covering 70% of the area under the WINE eigenvalues.

Scree-Plot Gap (SP-G). The criterion uses a greedy search for finding the largest variance difference between two consecutive pairs of scale-normalized eigenvalues [Zhu and Ghodsi 2006]. The intersection of both pairs is returned as d' . Figure 2(b) shows the search on WINE where line segments represent the scaled differences.

Plane visualization (2D). This approach is a baseline rule that sets $d' = 2$ so that data can be visualized in a Euclidean plane. Figure 1(a) shows that PCA reduction for set WINE.

Statistical criteria examine whether a set of eigenvalues is larger than an expected value drawn from a known data distribution. Approaches of that category include:

Kaiser-Guttman (KG). This rule retrieves every eigenvalue greater than 1.0 aiming at retaining shared data variance [Guttman 1954, Neto et al. 2005].

Broken-Stick (B-St). If joint variance is randomly distributed within d axes, then eigenvalues are supposed to follow the *Broken-Stick distribution* [Legendre and Legendre 2012]. For a set of eigenvalues e_i , Broken-Stick entries are given by $bs_i = \left(\frac{\sum_{j=1}^d e_j}{d} \right) / \left(\sum_{j=1}^{d-i+1} \frac{1}{d+1-j} \right)$. *B-St* returns the last scaled eigenvalue that deviates from the random distribution, *i.e.*, $d' = k \mid \left(e_{k+1} / \sum_{i=1}^d e_i \leq bs_{k+1} \wedge bs_1 \leq e_1 / \sum_{i=1}^d e_i \right) \vee \left(e_{k+1} / \sum_{i=1}^d e_i \geq bs_{k+1} \wedge bs_1 \geq e_1 / \sum_{i=1}^d e_i \right)$. Figure 2(c) depicts the *B-St* rule for set WINE.

Finally, intrinsic dimension-based criteria estimate possible correlations embedded within data features. A stable rule for obtaining such a value is as follows.

Rho-Score (ρ -sct). This estimator approximates data intrinsic dimension by using its distance distribution [Pestov 2008]. Formally, let $\mathcal{R} \subset \mathbb{R}^d$ be a dataset and $\delta : \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}_+$ be the Euclidean distance, the ρ -sct criterion estimates d' as in Eq. 1.

$$d' = \left\lceil \frac{\mu_{\mathcal{R}}^2}{2\sigma_{\mathcal{R}}^2} \right\rceil ; \mu_{\mathcal{R}} = \frac{1}{2|\mathcal{R}|} \sum_{r_i, r_j \in \mathcal{R}} \delta(r_i, r_j); \sigma_{\mathcal{R}} = \sqrt{\frac{1}{|\mathcal{R}| - 1} \sum_{r_i, r_j \in \mathcal{R}} \left(\frac{\delta(r_i, r_j) - \mu_{\mathcal{R}}}{2} \right)^2} \quad (1)$$

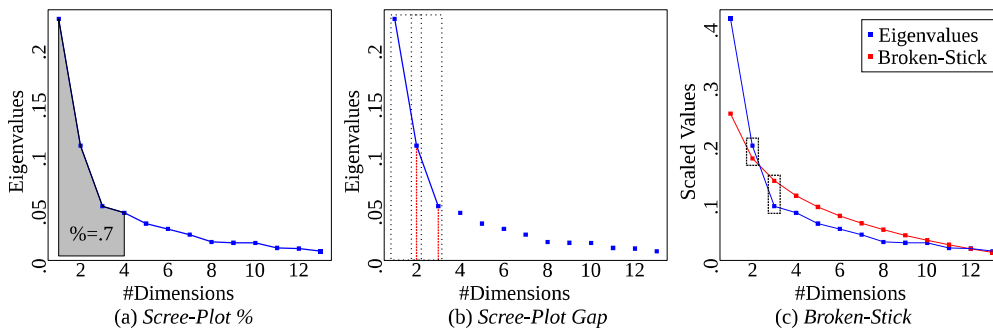


Figure 2. Choosing the number of dimensions d' in a PCA reduction of WINE.

3. Material and Methods

Estimation criteria implementation. We implement PCA, and the six reviewed estimation criteria from scratch in \mathbb{R}^2 by taking advantage of the spectral decomposition *eigen* routine available at R-core package. Graphical-based criteria are constructed upon sorted eigenvalues produced by the *eigen* routine, as well as the Kaiser-Guttman rule. An equi-width histogram H with d positions is employed as the underlying data structure for the *B-St* implementation. Each histogram position corresponds to a theoretical descending-sorted eigenvalue that follows a scale-normalized Broken-Stick distribution. Accordingly, the *B-St* estimation is carried out by a linear search over H for finding the first position that deviates from the eigenvalues' distribution. Finally, we implemented the Euclidean distance, as well as both mean and variance Welford's one-pass calculation for efficiently obtaining the intrinsic dimension returned by the ρ -*sct* rule.

Non-parametric and pairwise hypothesis tests. We adopt pairwise tests for assessing which criteria are suitable for estimating the number of dimensions in PCA reductions of labeled sets. The two-tailed Wilcoxon test is convenient for such a pairwise comparison because the null hypothesis is that results produced by two competing estimators are drawn from the same distribution, whereas the alternative hypothesis is that outputs are from different distributions [Wilcoxon 1992]. We model the results of PCA criteria as the quality values measured after a classifier 10-folds cross-validation, *e.g.*, F-Measure [Aggarwal 2015]. Therefore, a pair of criteria is compared by its paired list of label-driven results, being that pairwise list sorted by the absolute differences between the two juxtaposed outputs. An incremental *rank* ranging from one to the number of observations is assigned to each position of the sorted list. Ranks of positions where the second rule outperforms the first are multiplied by -1 , which generates two rank groups aggregated by a *ranking sum*. Such a sum adjusted by the number of observations produces a z -score, which is employed in the analysis of Wilcoxon's null hypothesis.

4. Experiments

We evaluate the performance of criteria *SP-%*, *SP-G*, *2D*, *KG*, *B-St*, and ρ -*sct* in association with four different classifiers NB, IbL, DT, and MLP in 20 medium to high dimensional labeled datasets (\mathcal{R}) from UCI, MILD, GBDI and QTDU repositories². In particular, we relied on both Weka v3.8.4 and R v3.6.1 to set up the classifiers as follows: (i) DT with binary splits and early pruning, (ii) MLP with one hidden layer of \sqrt{d} neurons, and (iii) IbL with Euclidean distance and one neighbor. Evaluations were conducted in a KUbuntu machine 19.10 with an Intel *i5* processor, 8GB RAM, and a 1TB disk.

We tuned the *SP-%* criterion by evaluating its performance with parameters $\% = \{.65, .7, .75\}$, being the setup $SP-\% = .7$ the tuning with highest F-Measures, on average. Table 1 details the F-Measure reached by each criterion for every evaluation scenario, where *Emb.* lines stand for data original dimensionality d . Results reinforce no monotonic relationship can be drawn from F-Measure and the number of dimensions (d').

Next, we compare the competing criteria through a set of pairwise Wilcoxon's tests grouped by classifiers. Figure 3 shows the z -scores obtained for each pairwise comparison (**line vs. column**), and highlights the cases in which the null hypothesis was re-

least 90% regarding NB and MLP classification, whereas $SP\%$ also outperformed estimators $SP-G$, KG , and $2D$. A similar result was observed for classifiers IbL and DT in which $B-St$ dominated $SP-G$, KG , $\rho-sct$, and $2D$ within significance levels. In those scenarios, $SP\%$ also outperformed $SP-G$, KG , $2D$, and $\rho-sct$ with statistical significance. Such findings pinpoint *both criteria $B-St$ and $SP\%$ are suitable rules* for choosing the number of dimensions in a PCA reduction. Notice, however, *$B-St$ is a parameterless estimator that may be preferable to $SP\%$ whenever adjusting variance area $\%$ is unpractical.*

A counterpart discovery is estimators KG and $SP-G$ performed poorly in comparison to other criteria regarding labeling-driven tasks. In particular, KG did not outperform any competitor, including $SP-G$, and baseline rules $2D$ and $Emb.$. Lastly, $\rho-sct$ criterion showed an intermediary performance, which indicates there may be a relationship between the intrinsic dimension and the number of PCA axes, but not strong as a correlation to be spotted by the experimental supervised evaluation we carried out.

5. Conclusions and Future Work

This study has discussed global criteria for finding the number of dimensions in a PCA reduction of labeled datasets. Since estimators are based on distinct theoretical grounds, we examine their performance from an experimental perspective regarding the biases of different classifiers. Results indicate $B-St$ and $SP\%$ are suitable rules for estimating the number of PCA axes, whereas KG and $SP-G$ shall be avoided in the reductions. Such outcomes enable devising the tuning of PCA-based pieces of AutoMLs in future work.

Acknowledgments. The study was supported by FAPERJ and CEPID-CeMEAI/FAPESP (Grants 2013/07375-0 – 2019/01735-0).

References

- Aggarwal, C. (2015). *Data mining: The textbook*. Springer.
- Guttman, L. (1954). Some necessary conditions for common-factor analysis. *Psychometrika*, 19(2):149–161.
- Jackson, D. A. (1993). Stopping rules in principal components analysis: A comparison of heuristical and statistical approaches. *Ecology*, 74(8):2204–2214.
- James, G., Witten, D., Hastie, T., and Tibshirani, R. (2013). *An introduction to Statistical Learning*, volume 112. Springer.
- Legendre, P. and Legendre, L. F. (2012). *Numerical Ecology*. Elsevier.
- Neto, P., Jackson, D., and Somers, K. (2005). How many principal components? Stopping rules for determining the number of non-trivial axes revisited. *C. Stat.*, 49(4):974–997.
- Pearson, K. (1901). On lines and planes of closest fit to systems of points in space. *The London, Edinburgh, and Dublin Phil. Magazine and J. of Science*, 2(11):559–572.
- Pestov, V. (2008). An axiomatic approach to intrinsic dimension of a dataset. *Neural Networks*, 21(2-3):204–213.
- Wilcoxon, F. (1992). Individual comparisons by ranking methods. In *Breakthroughs in statistics*, pages 196–202. Springer.
- Zhu, M. and Ghodsi, A. (2006). Automatic dimensionality selection from the scree plot via the use of profile likelihood. *Computational Stat. & Data Analysis*, 51(2):918–930.