


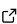

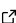
1 ExpFamilyPCA.jl: A Julia Package for Exponential 2 Family Principal Component Analysis

3 Logan Mondal Bhamidipaty ¹, Mykel J. Kochenderfer ¹, and Trevor
4 Hastie ¹

5 ¹ Stanford University

DOI: [10.xxxxxx/draft](https://doi.org/10.xxxxxx/draft)

Software

- [Review](#) 
- [Repository](#) 
- [Archive](#) 

Editor: 

Submitted: 12 October 2024

Published: unpublished

License

Authors of papers retain copyright
and release the work under a
Creative Commons Attribution 4.0
International License ([CC BY 4.0](https://creativecommons.org/licenses/by/4.0/)).

6 Summary

7 Principal component analysis (PCA) ([Hotelling, 1933](#); [Jolliffe, 2002](#); [Pearson, 1901](#)) is popular
8 for compressing, denoising, and interpreting high-dimensional data, but it underperforms on
9 binary, count, and compositional data because the objective assumes data is normally distributed.
10 Exponential family PCA (EPCA) ([Collins et al., 2001](#)) generalizes PCA to accommodate data
11 from any exponential family distribution, making it more suitable for fields where these data
12 types are common, such as geochemistry, marketing, genomics, political science, and machine
13 learning ([Greenacre, 2021](#); [Hastie et al., 2009](#)).

14 ExpFamilyPCA.jl is a library for EPCA written in Julia, a dynamic language for scientific
15 computing ([Bezanson et al., 2012](#)). It is the first EPCA package in Julia and the first in any
16 language to support EPCA for multiple distributions.

17 Statement of Need

18 EPCA is used in reinforcement learning ([Roy et al., 2005](#)), sample debiasing ([R. Huang &
19 Lee, 2023](#)), and compositional analysis ([Gan & Valdez, 2024](#)). Wider adoption, however,
20 remains limited due to the lack of implementations. The only other EPCA package is written
21 in MATLAB and supports just one distribution ([Chambrier, 2016](#)). This is surprising, as other
22 Bregman-based optimization techniques have been successful in areas like mass spectrometry
23 ([Nozaki & Nakamoto, 2017](#)), ultrasound denoising ([J. Huang & Yang, 2013](#)), topological
24 data analysis ([Edelsbrunner & Wagner, 2019](#)), and robust clustering ([Banerjee et al., 2005](#)).
25 These successes suggest that EPCA holds untapped potential in signal processing and machine
26 learning.

27 The absence of a general EPCA library likely stems from the limited interoperability between
28 fast symbolic differentiation and optimization libraries in popular languages like Python and
29 C. Julia, by contrast, uses multiple dispatch which promotes high levels of generic code
30 reuse ([Karpinski, 2019](#)). Multiple dispatch allows ExpFamilyPCA.jl to integrate fast symbolic
31 differentiation ([Gowda et al., 2022](#)), optimization ([Mogensen & Riseth, 2018](#)), and numerically
32 stable computation ([Mächler, 2015](#)) without requiring costly API conversions. As a result,
33 ExpFamilyPCA.jl delivers speed, stability, and flexibility, with built-in support for most common
34 distributions (§ [Supported Distributions](#)) and flexible constructors for custom distributions (§
35 [Custom Distributions](#)).

36 **Principal Component Analysis**

37 **Geometric Interpretation**

38 Given a data matrix $X \in \mathbb{R}^{n \times d}$ with n observations and d features, PCA seeks the closest
39 low-rank approximation $\Theta \in \mathbb{R}^{n \times d}$ by minimizing the reconstruction error

$$\begin{aligned} & \underset{\Theta}{\text{minimize}} && \frac{1}{2} \|X - \Theta\|_F^2 \\ & \text{subject to} && \text{rank}(\Theta) = k \end{aligned}$$

40 where $\|\cdot\|_F$ denotes the Frobenius norm. The optimal Θ is a k -dimensional linear subspace
41 that can be written as the product of the projected observations $A \in \mathbb{R}^{n \times k}$ and the basis
42 $V \in \mathbb{R}^{k \times d}$:

$$X \approx \Theta = AV.$$

43 This suggests that each observation $x_i \in \text{rows}(X)$ can be well-approximated by a linear
44 combination of k basis vectors (the rows of V):

$$x_i \approx \theta_i = a_i V$$

45 for $i = 1, \dots, n$.

46 **Probabilistic Interpretation**

47 The PCA objective is equivalent to maximum likelihood estimation for a Gaussian model.
48 Under this lens, each observation x_i is a noisy realization of a d -dimensional Gaussian at
49 $\theta_i \in \text{rows}(\Theta)$:

$$x_i \sim \mathcal{N}(\theta_i, I).$$

50 To recover the latent structure Θ , PCA solves

$$\begin{aligned} & \underset{\Theta}{\text{maximize}} && \sum_{i=1}^n \log \mathcal{L}(x_i; \theta_i) \\ & \text{subject to} && \text{rank}(\Theta) = k \end{aligned}$$

51 where \mathcal{L} is the likelihood function.

52 **Exponential Family PCA**

53 **Link Function**

54 The link function $g(\theta)$ connects the natural parameter θ to the mean parameter μ of an
55 exponential family distribution. It is defined as the gradient of the log-partition function $G(\theta)$:

$$\mu = g(\theta) = \nabla G(\theta).$$

56 The link function serves a role analogous to that in generalized linear models (GLMs) (McCullagh
57 & Nelder, 1989). In GLMs, the link function connects the linear predictor to the mean of
58 the distribution, enabling flexibility in modeling various data types. Similarly, in EPCA, the
59 link function maps the low-dimensional latent variables to the expectation parameters of
60 the exponential family, thereby generalizing the linear assumptions of traditional PCA to
61 accommodate diverse distributions (see [appendix](#)).

62 **Bregman Divergences**

63 EPCA extends the probabilistic interpretation of PCA using a measure of statistical difference
 64 called the Bregman divergence (Bregman, 1967; Efron, 2004). The Bregman divergence B_F
 65 for a strictly convex, continuously differentiable function F is

$$B_F(p\|q) = F(p) - F(q) - \langle \nabla F(q), p - q \rangle.$$

66 This can be interpreted as the difference between $F(p)$ and its linear approximation about
 67 q . When F is the convex conjugate of the log-partition function of an exponential family
 68 distribution, minimizing the Bregman divergence corresponds to maximizing the associated
 69 log-likelihood (Azoury & Warmuth, 2001; Forster & Warmuth, 2002) (see [documentation](#)).

70 **Loss Function**

71 EPCA generalizes the PCA objective as a Bregman divergence between the data X and the
 72 expectation parameters $g(\Theta)$:

$$\begin{aligned} & \underset{\Theta}{\text{minimize}} && B_F(X\|g(\Theta)) \\ & \text{subject to} && \text{rank}(\Theta) = k \end{aligned}$$

73 where

- 74 ■ $g(\theta)$ is the **link function** and the gradient of G ,
- 75 ■ $G(\theta)$ is a strictly convex, continuously differentiable function (usually the **log-partition**
 76 of an exponential family distribution),
- 77 ■ and $F(\mu)$ is the **convex conjugate** of G defined by

$$F(\mu) = \langle \mu, \theta \rangle - G(\theta).$$

78 This suggests that data from the exponential family is well-approximated by expectation
 79 parameters

$$x_i \approx g(\theta_i) = g(a_i V).$$

80 **Regularization**

81 To ensure the optimum converges, we introduce a regularization term

$$\begin{aligned} & \underset{\Theta}{\text{minimize}} && B_F(X\|g(\Theta)) + \epsilon B_F(\mu_0\|g(\Theta)) \\ & \text{subject to} && \text{rank}(\Theta) = k \end{aligned}$$

82 where $\epsilon > 0$ and $\mu_0 \in \text{range}(g)$.

83 **Example: Poisson EPCA**

84 The Poisson EPCA objective is the generalized Kullback-Leibler (KL) divergence (see [appendix](#)),
 85 making Poisson EPCA ideal for compressing discrete distribution data.

86 This is useful in applications like belief compression in reinforcement learning (Roy et al.,
 87 2005), where high-dimensional belief states can be effectively reduced with minimal information
 88 loss. Below we recreate a figure from Roy & Gordon (2002) and observe that Poisson EPCA
 89 achieved a nearly perfect reconstruction of a 41-dimensional belief profile using just 5 basis
 90 components.

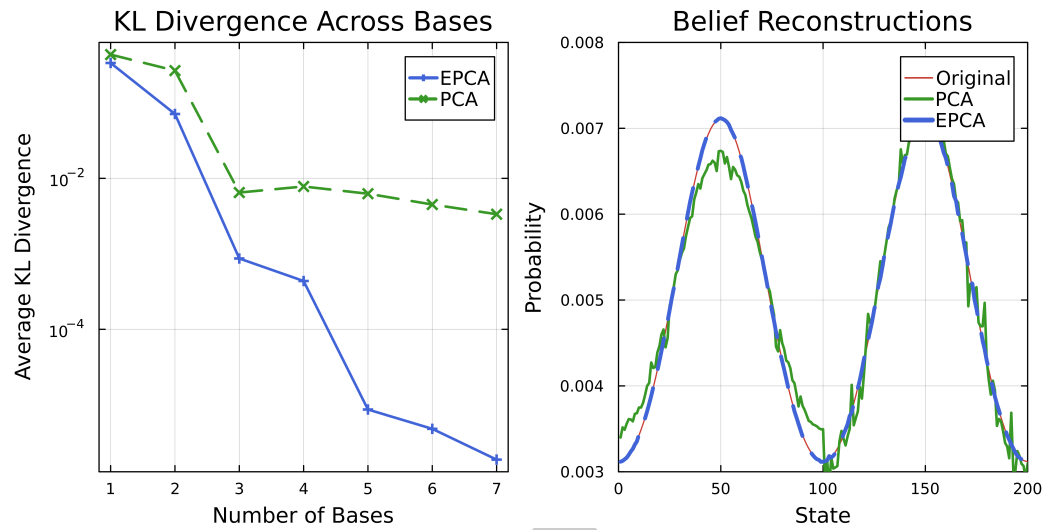


Figure 1: Left - KL Divergence for Poisson EPCA versus PCA. Right - Reconstructions from the models.

91 For a larger environment with 200 states, PCA struggles even with 10 basis.

92 API

93 Supported Distributions

94 ExpFamilyPCA.jl includes efficient EPCA implementations for several exponential family
95 distributions.

Julia	Description
BernoulliEPCA	For binary data
BinomialEPCA	For count data with a fixed number of trials
ContinuousBernoulliEPCA	For modeling probabilities between 0 and 1
GammaEPCA	For positive continuous data
GaussianEPCA	Standard PCA for real-valued data
NegativeBinomialEPCA	For over-dispersed count data
ParetoEPCA	For modeling heavy-tailed distributions
PoissonEPCA	For count and discrete distribution data
WeibullEPCA	For modeling life data and survival analysis

96 Custom Distributions

97 When working with custom distributions, certain specifications are often more convenient and
98 computationally efficient than others. For example, inducing the gamma EPCA objective from
99 the log-partition $G(\theta) = -\log(-\theta)$ and its derivative $g(\theta) = -1/\theta$ is much simpler than
100 implementing the full the Itakura-Saito distance (Itakura & Saito, 1968) (see [appendix](#)):

$$D(P(\omega), \hat{P}(\omega)) = \frac{1}{2\pi} \int_{-\pi}^{\pi} \left[\frac{P(\omega)}{\hat{P}(\omega)} - \log \frac{P(\omega)}{\hat{P}(\omega)} - 1 \right] d\omega.$$

101 In ExpFamilyPCA.jl, we would write:

```
G(θ) = -log(-θ)
g(θ) = -1 / θ
gamma_epca = EPCA(indim, outdim, G, g, Val(:(G, :g)); options = NegativeDomain())
```

102 A lengthier discussion of the EPCA constructors and math is provided in the [documentation](#).

103 Usage

104 Each EPCA object supports a three-method interface: `fit!`, `compress`, and `decompress`. `fit!`
105 trains the model and returns the compressed training data; `compress` returns compressed input;
106 and `decompress` reconstructs the original data from the compressed representation.

```
X = sample_from_gamma(n1, indim)
```

```
Y = sample_from_gamma(n2, indim)
```

```
X_compressed = fit!(gamma_epca, X)
```

```
Y_compressed = compress(gamma_epca, Y)
```

```
Y_reconstructed = decompress(gamma_epca, Y_compressed)
```

107 Acknowledgments

108 We thank Ryan Tibshirani, Alec Jamgochian, Robert Moss, and Dylan Asmar for their help
109 and guidance.

110 References

111 Azoury, K. S., & Warmuth, M. K. (2001). Relative loss bounds for on-line density esti-
112 mation with the exponential family of distributions. *Machine Learning*, 43, 211–246.
113 <https://doi.org/https://doi.org/10.1023/A:1010896012157>

114 Banerjee, A., Merugu, S., Dhillon, I. S., & Ghosh, J. (2005). Clustering with Bregman
115 divergences. *Journal of Machine Learning Research*, 6(58), 1705–1749. [http://jmlr.org/
116 papers/v6/banerjee05b.html](http://jmlr.org/papers/v6/banerjee05b.html)

117 Bezanson, J., Karpinski, S., Shah, V. B., & Edelman, A. (2012). *Julia: A fast dynamic*
118 *language for technical computing*. <https://doi.org/10.48550/arXiv.1209.5145>

119 Bregman, L. M. (1967). The relaxation method of finding the common point of convex
120 sets and its application to the solution of problems in convex programming. *USSR*
121 *Computational Mathematics and Mathematical Physics*, 7(3), 200–217. [https://doi.org/
122 10.1016/0041-5553\(67\)90040-7](https://doi.org/10.1016/0041-5553(67)90040-7)

123 Chambrier, G. de. (2016). *E-PCA*. <https://github.com/gpldecha/e-pca>

124 Collins, M., Dasgupta, S., & Schapire, R. E. (2001). A generalization of principal components
125 analysis to the exponential family. *Advances in Neural Information Processing Systems*, 14.

126 Edelsbrunner, H., & Wagner, H. (2019). Topological data analysis with Bregman divergences.
127 *Journal of Computational Geometry*, Vol. 9 No. 2 (2018): Special Issue of Selected Papers
128 from SoCG 2017. <https://doi.org/10.20382/JOCG.V9I2A6>

129 Efron, B. (2004). The estimation of prediction error. *Journal of the American Statistical*
130 *Association*, 99(467), 619–632. <https://doi.org/10.1198/016214504000000692>

131 Forster, J., & Warmuth, M. K. (2002). Relative expected instantaneous loss bounds. *Journal*
132 *of Computer and System Sciences*, 64(1), 76–102. <https://doi.org/10.1006/jcss.2001.1798>

133 Gan, G., & Valdez, E. A. (2024). Compositional Data Regression in Insurance with Exponential
134 Family PCA. *Variance*, 17(1).

- 135 Gowda, S., Ma, Y., Cheli, A., Gwóźdzdź, M., Shah, V. B., Edelman, A., & Rackauckas,
136 C. (2022). High-performance symbolic-numeric via multiple dispatch. *Association for*
137 *Computing Machinery Communications in Computer Algebra*, 55(3), 92–96. <https://doi.org/10.1145/3511528.3511535>
138
- 139 Greenacre, M. (2021). Compositional data analysis. *Annual Review of Statistics and Its*
140 *Application*, 8(1), 271–299.
- 141 Hastie, T., Tibshirani, R., Friedman, J. H., & Friedman, J. H. (2009). *The Elements of*
142 *Statistical Learning: Data Mining, Inference, and Prediction* (Vol. 2). Springer.
- 143 Hotelling, H. (1933). Analysis of a complex of statistical variables into principal components.
144 *Journal of Educational Psychology*, 24, 498–520. [https://api.semanticscholar.org/CorpusID:](https://api.semanticscholar.org/CorpusID:144828484)
145 [144828484](https://api.semanticscholar.org/CorpusID:144828484)
- 146 Huang, J., & Yang, X. (2013). Fast reduction of speckle noise in real ultrasound images.
147 *Signal Processing*, 93(4), 684–694. <https://doi.org/10.1016/j.sigpro.2012.09.005>
- 148 Huang, R., & Lee, Y. (2023). *Debiasing sample loadings and scores in exponential family PCA*
149 *for sparse count data*. <https://arxiv.org/abs/2312.13430>
- 150 Itakura, F., & Saito, S. (1968). Analysis synthesis telephony based on the maximum likelihood
151 method. *Proceedings of the 6th International Congress on Acoustics*, C-17-C-20.
- 152 Jolliffe, I. T. (2002). *Principal component analysis for special types of data*. Springer.
- 153 Karpinski, S. (2019). *The unreasonable effectiveness of multiple dispatch*. Conference Talk at
154 JuliaCon 2019, available at <https://www.youtube.com/watch?v=kc9HwsxE1OY>.
- 155 Mächler, M. (2015). *Accurately computing $\log(1 - \exp(-|a|))$ assessed by the 'Rmpfr' package*.
156 <https://doi.org/10.13140/RG.2.2.11834.70084>
- 157 McCullagh, P., & Nelder, J. A. (1989). *Generalized Linear Models* (2nd ed.). Chapman &
158 Hall/CRC.
- 159 Mogensen, P. K., & Riseth, A. N. (2018). Optim: A mathematical optimization package for
160 Julia. *Journal of Open Source Software*, 3(24), 615. <https://doi.org/10.21105/joss.00615>
- 161 Nozaki, Y., & Nakamoto, T. (2017). Itakura-Saito distance based autoencoder for dimension-
162 ality reduction of mass spectra. *Chemometrics and Intelligent Laboratory Systems*, 167,
163 63–68. <https://doi.org/10.1016/j.chemolab.2017.05.002>
- 164 Pearson, K. (1901). On lines and planes of closest fit to systems of points in space. *The*
165 *London, Edinburgh, and Dublin Philosophical Magazine and Journal of Science*, 2(11),
166 559–572. <https://doi.org/10.1080/14786440109462720>
- 167 Roy, N., & Gordon, G. (2002). Exponential family PCA for belief compression in POMDPs. In
168 S. Becker, S. Thrun, & K. Obermayer (Eds.), *Advances in Neural Information Processing*
169 *Systems* (Vol. 15). MIT Press. [https://proceedings.neurips.cc/paper_files/paper/2002/](https://proceedings.neurips.cc/paper_files/paper/2002/file/a11f9e533f28593768ebf87075ab34f2-Paper.pdf)
170 [file/a11f9e533f28593768ebf87075ab34f2-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2002/file/a11f9e533f28593768ebf87075ab34f2-Paper.pdf)
- 171 Roy, N., Gordon, G., & Thrun, S. (2005). Finding approximate POMDP solutions through
172 belief compression. *Journal of Artificial Intelligence Research*, 23, 1–40. [https://doi.org/](https://doi.org/10.1613/jair.1496)
173 [10.1613/jair.1496](https://doi.org/10.1613/jair.1496)