

Notes on Archetypal Analysis

Philipp Sven Lars Schäfer

April 2025

1 Problem Formulation

Let $\mathcal{X} = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}_{n=1}^N$ be a data set consisting of N D -dimensional data points, and let $\mathbf{X} \in \mathbb{R}^{N \times D}$ be the matrix where each row is a data point.

In Archetypal Analysis we make two assumptions:

1. Each data point is a convex combination of K archetypes;
2. Each archetype is a convex combination of N data points.

Expressing the first assumption in matrix notation yields

$$\hat{\mathbf{X}} = \mathbf{AZ} \quad \text{or} \quad \hat{\mathbf{x}}_n = \mathbf{Z}^T \mathbf{a}_n \quad \text{for } n = 1, \dots, N \quad (1)$$

where $\hat{\mathbf{X}} \in \mathbb{R}^{N \times D}$ is the reconstructed data matrix, $\mathbf{Z} \in \mathbb{R}^{K \times D}$ is the matrix of archetypes (i.e. each row is one archetype), and $\mathbf{A} \in \mathbb{R}^{N \times K}$ is a row-stochastic matrix that defines by which archetypes each data point is formed.

Expressing the second assumption in matrix notation yields

$$\mathbf{Z} = \mathbf{BX} \quad \text{or} \quad \mathbf{z}_k = \mathbf{X}^T \mathbf{b}_k \quad \text{for } k = 1, \dots, K \quad (2)$$

where $\mathbf{B} \in \mathbb{R}^{K \times N}$ is a row-stochastic matrix that defines by which data point each archetype is defined by.

The reconstruction error is most commonly measured using the residual sum of squares (RSS), given by the squared Frobenius norm,

$$\|\mathbf{X} - \hat{\mathbf{X}}\|_F^2 = \|\mathbf{X} - \mathbf{AZ}\|_F^2 = \|\mathbf{X} - \mathbf{ABX}\|_F^2 \quad (3)$$

which yields the following optimization objective

$$\begin{aligned} \mathbf{A}^*, \mathbf{B}^* = \arg \min_{\substack{\mathbf{A} \in \mathbb{R}^{N \times K} \\ \mathbf{B} \in \mathbb{R}^{K \times N}}} \|\mathbf{X} - \mathbf{ABX}\|_F^2 \quad \text{subject to} \\ \mathbf{A} \geq 0, \mathbf{A}\mathbf{1}_K = \mathbf{1}_N \\ \mathbf{B} \geq 0, \mathbf{B}\mathbf{1}_N = \mathbf{1}_K \end{aligned} \quad (4)$$

Introducing the set of row-stochastic non-negative matrices,

$$F(N, K) := \{\mathbf{A} \in \mathbb{R}^{N \times K} \mid \mathbf{A} \geq 0 \wedge \mathbf{A}\mathbf{1}_K = \mathbf{1}_N\} \quad (5)$$

we can write the objective compactly as:

$$\mathbf{A}^*, \mathbf{B}^* = \arg \min_{\substack{\mathbf{A} \in F(N, K) \\ \mathbf{B} \in F(K, N)}} \|\mathbf{X} - \mathbf{ABX}\|_F^2 \quad (6)$$

2 Properties of the Objective

Property 1 (Translation invariance): The minimizers $\mathbf{A}^*, \mathbf{B}^*$ of the objective are invariant under row-wise translations of \mathbf{X} . Let $\tilde{\mathbf{X}} = \mathbf{X} + \mathbf{1}_N \mathbf{v}^T$ for any $\mathbf{v} \in \mathbb{R}^D$, then

$$\arg \min_{\substack{\mathbf{A} \in F(N, K) \\ \mathbf{B} \in F(K, N)}} \|\tilde{\mathbf{X}} - \mathbf{AB}\tilde{\mathbf{X}}\|_F^2 = \arg \min_{\substack{\mathbf{A} \in F(N, K) \\ \mathbf{B} \in F(K, N)}} \|\mathbf{X} - \mathbf{ABX}\|_F^2 \quad (7)$$

Proof: Let $\mathbf{v} \in \mathbb{R}^D$, and let $\tilde{\mathbf{X}} = \mathbf{X} + \mathbf{1}_N \mathbf{v}^T$ be the translated matrix. Then for any feasible \mathbf{A}, \mathbf{B}

$$\begin{aligned} \tilde{\mathbf{X}} - \mathbf{AB}\tilde{\mathbf{X}} &= (\mathbf{X} + \mathbf{1}_N \mathbf{v}^T) - \mathbf{AB}(\mathbf{X} + \mathbf{1}_N \mathbf{v}^T) \\ &= \mathbf{X} + \mathbf{1}_N \mathbf{v}^T - \mathbf{ABX} - \mathbf{AB}\mathbf{1}_N \mathbf{v}^T \end{aligned} \quad (8)$$

Since $\mathbf{B}\mathbf{1}_N = \mathbf{1}_K$ and $\mathbf{A}\mathbf{1}_K = \mathbf{1}_N$, this simplifies to

$$\begin{aligned} \tilde{\mathbf{X}} - \mathbf{AB}\tilde{\mathbf{X}} &= \mathbf{X} + \mathbf{1}_N \mathbf{v}^T - \mathbf{ABX} - \mathbf{1}_N \mathbf{v}^T \\ &= \mathbf{X} - \mathbf{ABX} \end{aligned} \quad (9)$$

Therefore, the reconstruction error remains unchanged, and the minimizers $\mathbf{A}^*, \mathbf{B}^*$ are invariant under such translations. Thus, the minimizers $\mathbf{A}^*, \mathbf{B}^*$ are invariant to centering the data.

Property 2 (Scale invariance): The minimizers $\mathbf{A}^*, \mathbf{B}^*$ of the objective are invariant under global scaling of \mathbf{X} . Let $\tilde{\mathbf{X}} = \lambda \mathbf{X}$ for any $\lambda \neq 0$, then

$$\arg \min_{\substack{\mathbf{A} \in F(N, K) \\ \mathbf{B} \in F(K, N)}} \|\tilde{\mathbf{X}} - \mathbf{AB}\tilde{\mathbf{X}}\|_F^2 = \arg \min_{\substack{\mathbf{A} \in F(N, K) \\ \mathbf{B} \in F(K, N)}} \|\mathbf{X} - \mathbf{ABX}\|_F^2 \quad (10)$$

Proof: Let $\lambda \neq 0$, and let $\tilde{\mathbf{X}} = \lambda \mathbf{X}$ be the scaled matrix. Then for any feasible \mathbf{A}, \mathbf{B}

$$\begin{aligned} \tilde{\mathbf{X}} - \mathbf{AB}\tilde{\mathbf{X}} &= \lambda \mathbf{X} - \mathbf{AB}\lambda \mathbf{X} \\ &= \lambda (\mathbf{X} - \mathbf{ABX}) \end{aligned} \quad (11)$$

Thus the objective for the scaled matrix is given by

$$\arg \min_{\substack{\mathbf{A} \in F(N, K) \\ \mathbf{B} \in F(K, N)}} \|\tilde{\mathbf{X}} - \mathbf{AB}\tilde{\mathbf{X}}\|_F^2 = \arg \min_{\substack{\mathbf{A} \in F(N, K) \\ \mathbf{B} \in F(K, N)}} \lambda^2 \|\mathbf{X} - \mathbf{ABX}\|_F^2 \quad (12)$$

Since $\lambda \neq 0$, we have $\lambda^2 > 0$, and thus the objective is scaled by a positive constant. Multiplying the objective function by a positive scalar does not affect the location of its minimum, because the ordering of objective values is preserved. In particular, the first-order (stationarity) and second-order (convexity/curvature) necessary conditions for optimality remain unchanged under such scaling.

Property 3 (unique up to permutation of archetypes / No rotational ambiguity)

Assuming that for each archetypes there exists one data point that only belong to this archetype

$$\forall k \in \{1, \dots, K\} \exists n \in \{1, \dots, N\} a_{nk} > 0 \wedge a_{ak'} = 0 \forall k' \neq k \quad (13)$$

and that for each archetype there exists one data point that is only used to define this archetype and not any other archetype

$$\forall k \in \{1, \dots, K\} \exists n \in \{1, \dots, N\} z_{kn} > 0 \wedge z_{k'n} = 0 \forall k' \neq k \quad (14)$$

then the objective does not suffer from rotational ambiguity.

Note, these conditions mean that both \mathbf{A} and \mathbf{B} have rank K (i.e. K linearly indendent columns / rows).

Let $\mathbf{Q} \in \mathbb{R}^{K \times K}$ be some invertible matrix

$$\mathbf{A}\mathbf{B}\mathbf{X} = \mathbf{A}\mathbf{Q}\mathbf{Q}^{-1}\mathbf{B}\mathbf{X} = \tilde{\mathbf{A}}\tilde{\mathbf{B}}\mathbf{X} \quad (15)$$

Requiring that $\tilde{\mathbf{A}} \in F(N, K)$ and $\tilde{\mathbf{B}} \in F(K, N)$ (i.e. that both $\tilde{\mathbf{A}}, \tilde{\mathbf{B}}$ are still row-stochastic), we can derive the following properties that \mathbf{Q} must fullfull.

First, since we require $\tilde{\mathbf{A}} \geq 0$, and $\mathbf{A} \geq 0$, and $\tilde{\mathbf{A}} = \mathbf{A}\mathbf{Q}^{-1}$, and Equation (13) must hold, we know that $\mathbf{Q} \geq 0$.

Second, since we require $\tilde{\mathbf{B}} \geq 0$, and $\mathbf{B} \geq 0$, and $\tilde{\mathbf{B}} = \mathbf{Q}^{-1}\mathbf{B}$, and Equation (14) must hold, we know that $\mathbf{Q}^{-1} \geq 0$.

Then, since \mathbf{Q} and \mathbf{Q}^{-1} are non-negative, Lemma 1.1 from [6] states that \mathbf{Q} must be a generalized permutation matrix, i.e. there exists some diagonal matrix $\mathbf{D} \in \mathbb{R}^{K \times K}$ and permutation matrix $\mathbf{P} \in \mathbb{R}^{K \times K}$ such that $\mathbf{Q} = \mathbf{D}\mathbf{P}$.

Third, since we require $\tilde{\mathbf{A}}\mathbf{1}_K = \mathbf{1}_N = \mathbf{A}\mathbf{Q}\mathbf{1}_K = \mathbf{1}_N$, we know that $\mathbf{Q}\mathbf{1}_K = \mathbf{1}_K$ and thus $\mathbf{Q} \in F(K, K)$ (i.e. \mathbf{Q} must be a row-stochastic matrix)

Then, this means that \mathbf{Q} must be a permutation matrix since

$$\begin{aligned} \mathbf{Q}\mathbf{1}_K &= \mathbf{1}_K \\ \rightarrow \mathbf{D}\mathbf{P}\mathbf{1}_K &= \mathbf{1}_K \\ \rightarrow \mathbf{D}\mathbf{1}_K &= \mathbf{1}_K \\ \rightarrow \mathbf{D} &= \mathbf{I}_K \end{aligned} \quad (16)$$

Property 4 (Rewrite via convex hull of \mathbf{Z}): For some fixed $\mathbf{B} \in F(K, N)$, let the corresponding archetype matrix be $\mathbf{Z} = \mathbf{B}\mathbf{X}$. Then for the optimal $\mathbf{A} \in F(K, N)$, the objective dfgafa

$$\begin{aligned}
\|\mathbf{X} - \mathbf{A}\mathbf{B}\mathbf{X}\|_F^2 &= \sum_{n=1}^N \|\mathbf{x}_n - \mathbf{Z}^T \mathbf{a}_n\|_2^2 \\
&= \sum_{n=1}^N \left\| \mathbf{x}_n - \sum_{k=1}^K a_{nk} \mathbf{z}_k \right\|_2^2
\end{aligned} \tag{17}$$

is equivalent to

$$\|\mathbf{X} - \mathbf{A}\mathbf{B}\mathbf{X}\|_F^2 = \sum_{n=1}^N \min_{\mathbf{q} \in \text{conv}(\mathcal{Z})} \|\mathbf{x}_n - \mathbf{q}\|_2^2 \tag{18}$$

Thus, for fixed \mathbf{B} , the optimal \mathbf{A} assigns each data point to its Euclidean projection onto the convex hull of the archetypes $\mathbf{z}_1, \dots, \mathbf{z}_K$. Any point $x_n \in \text{conv}(\mathcal{Z})$ does not contribute to the loss.

3 Optimization

While this objective is an Euclidean sum of square clustering problem which have been proven to be NP-hard[1], several practical optimization approaches have been developed that exploit that this objective is biconvex, meaning that it is convex in \mathbf{A} if we fix \mathbf{B} and vice versa. See Section 5 in Cutler & Breiman (1994) [3] or Section 2 in Mørup & Hansen (2012) [7] for more details. One way to optimize such a biconvex objective is to initialize \mathbf{A} , \mathbf{B} , and then alternating between solving the convex optimization problem in one variable fixing the other variable, and vice versa.

3.1 Gradient of the Objective

To compute the gradient of the unconstrained objective w.r.t. \mathbf{A} and \mathbf{B} , we first rewrite the residual sum of squares (Frobenius norm) in Equation (4) in terms of the trace

$$\begin{aligned}
\text{RSS} &= \|\mathbf{X} - \mathbf{A}\mathbf{B}\mathbf{X}\|_F^2 \\
&= \text{tr} \left((\mathbf{X} - \mathbf{A}\mathbf{B}\mathbf{X})^T (\mathbf{X} - \mathbf{A}\mathbf{B}\mathbf{X}) \right) \\
&= \text{tr}(\mathbf{X}^T \mathbf{X}) - \text{tr}(\mathbf{X}^T \mathbf{A}\mathbf{B}\mathbf{X}) - \text{tr}(\mathbf{X}^T \mathbf{B}^T \mathbf{A}^T \mathbf{X}) + \text{tr}(\mathbf{X}^T \mathbf{B}^T \mathbf{A}^T \mathbf{A}\mathbf{B}\mathbf{X}) \\
&= \text{tr}(\mathbf{X}^T \mathbf{X}) - 2 \text{tr}(\mathbf{X}^T \mathbf{A}\mathbf{B}\mathbf{X}) + \text{tr}(\mathbf{X}^T \mathbf{B}^T \mathbf{A}^T \mathbf{A}\mathbf{B}\mathbf{X})
\end{aligned} \tag{5}$$

where we used that for any $\mathbf{G}, \mathbf{H} \in \mathbb{R}^{N \times N}$ it is true that $\text{tr}(\mathbf{G} + \mathbf{H}) = \text{tr}(\mathbf{G}) + \text{tr}(\mathbf{H})$ and $\text{tr}(\mathbf{G}^T) = \text{tr}(\mathbf{G})$

Next we will use Equation 101 from the Matrix Cookbook by Petersen and Pedersen (2012) [8] which states that for any matrices $G, H, J \in \mathbb{R}^{N \times N}$ we have

$$\frac{\partial}{\partial H} \text{tr}(GHJ) = G^T J^T \tag{19}$$

and Equation 116 which states that for any matrices $G, H, J \in \mathbb{R}^{N \times N}$ we have

$$\frac{\partial}{\partial H} \text{tr}(G^T H^T JHG) = J^T HGG^T + JHGG^T \tag{20}$$

So computing the gradient of the RSS w.r.t. A we have

$$\begin{aligned}
G^{(A)} &= \nabla_A \text{RSS} \\
&= \nabla_A [\text{tr}(X^T X) - 2 \text{tr}(X^T A B X) + \text{tr}(X^T B^T A^T A B X)] \\
&= -2 \nabla_A \text{tr}(\underbrace{X^T}_G \underbrace{A}_H \underbrace{B X}_J) + \nabla_A \text{tr}(\underbrace{(B X)^T}_{G^T} \underbrace{A^T}_{H^T} \underbrace{I}_J \underbrace{A}_H \underbrace{B X}_G) \\
&= -2 X X^T B^T + (I^T A B X X^T B^T + I A B X X^T B^T) \\
&= -2 X X^T B^T + 2 A B X X^T B^T \\
&= 2 (A B X X^T B^T - X X^T B^T) \\
&= 2 (A Z Z^T - X Z^T)
\end{aligned} \tag{21}$$

Similarly, computing the gradient of the RSS w.r.t. B we have

$$\begin{aligned}
G^{(B)} &= \nabla_B \text{RSS} \\
&= \nabla_B [\text{tr}(X^T X) - 2 \text{tr}(X^T A B X) + \text{tr}(X^T B^T A^T A B X)] \\
&= -2 \nabla_B \text{tr}(\underbrace{X^T}_G \underbrace{A}_H \underbrace{B}_J \underbrace{X}_G) + \nabla_B \text{tr}(\underbrace{X^T}_{G^T} \underbrace{B^T}_{H^T} \underbrace{A^T}_J \underbrace{A}_H \underbrace{B}_G \underbrace{X}_G) \\
&= -2 A^T X X^T + (A^T A B X X^T + A^T A B X X^T) \\
&= -2 A^T X X^T + 2 A^T A B X X^T \\
&= 2 (A^T A B X X^T - A^T X X^T)
\end{aligned} \tag{22}$$

3.2 Regularized Nonnegative Least Squares

Introduced in 1994 by Adele Cutler and Leo Breiman [3], this was the first algorithm to solve the archetypal analysis objective in Equation (4).

The authors originally proposed to solve the constrained optimization problems using a Nonnegative Least Squares Problem (NNLS) solver and enforcing the convexity constraints using a penalty term with regularization parameter λ , i.e.

$$\begin{aligned}
\mathbf{a}_n &= \arg \min_{\mathbf{a}_n \in \mathbb{R}^K} \|\mathbf{x}_n - \mathbf{Z}^T \mathbf{a}_n\|_2^2 + \lambda \|\mathbf{1}_K - \mathbf{a}_n\|_2^2 \quad \text{subject to} \quad \mathbf{a}_n \geq 0 \\
&= \arg \min_{\mathbf{a}_n \in \mathbb{R}^K} \left\| \begin{bmatrix} \mathbf{x}_n \\ \lambda \end{bmatrix} - \begin{bmatrix} \mathbf{Z}^T \\ \lambda \mathbf{1}_K^T \end{bmatrix} \mathbf{a}_n \right\|_2^2
\end{aligned} \tag{23}$$

Equivalently, for \mathbf{B} we have

$$\begin{aligned}
\mathbf{b}_k &= \arg \min_{\mathbf{b}_k \in \mathbb{R}^N} \|\mathbf{z}_k - \mathbf{X}^T \mathbf{b}_k\|_2^2 + \lambda \|\mathbf{1}_N - \mathbf{b}_k\|_2^2 \quad \text{subject to} \quad \mathbf{b}_k \geq 0 \\
&= \arg \min_{\mathbf{b}_k \in \mathbb{R}^N} \left\| \begin{bmatrix} \mathbf{z}_k \\ \lambda \end{bmatrix} - \begin{bmatrix} \mathbf{X}^T \\ \lambda \mathbf{1}_N^T \end{bmatrix} \mathbf{b}_k \right\|_2^2
\end{aligned} \tag{24}$$

3.3 Principal Convex Hull Algorithm (PCHA)

Inspired by the projected gradient method for NMF [5] and normalization invariance approach introduced for NMF [4], the PCHA algorithm was introduced by Morten Mørup and Lars Kai Hansen in 2012 to solve the archetypal analysis objective.

Algorithm 1 Archetypal Analysis Algorithm

-
- 1: Initialize \mathbf{B} and compute the archetypes $\mathbf{Z} = \mathbf{B}\mathbf{X}$
 - 2: **while** not converged or maximum number of iterations is reached **do**
 - 3: **for** $n = 1$ to N **do**
 - 4: Find optimal \mathbf{a}_n by solving the constrained optimization problem:

$$\mathbf{a}_n = \arg \min_{\mathbf{a}_n \in \mathbb{R}^K} \|\mathbf{x}_n - \mathbf{Z}^T \mathbf{a}_n\|_2^2 \quad \text{subject to} \quad \mathbf{a}_n \geq 0, \sum_{k=1}^K a_{nk} = 1$$

- 5: **end for**
- 6: Compute the optimal archetypes \mathbf{Z} given \mathbf{A} , i.e.

$$\mathbf{Z} = \arg \min_{\mathbf{Z} \in \mathbb{R}^{K \times D}} \|\mathbf{X} - \mathbf{A}\mathbf{Z}\|_F^2$$

- 7: **for** $k = 1$ to K **do**
- 8: Find optimal \mathbf{b}_k by solving the constrained optimization problem:

$$\mathbf{b}_k = \arg \min_{\mathbf{b}_k \in \mathbb{R}^N} \|\mathbf{z}_k - \mathbf{X}^T \mathbf{b}_k\|_2^2 \quad \text{subject to} \quad \mathbf{b}_k \geq 0, \sum_{n=1}^N b_{kn} = 1$$

- 9: **end for**
 - 10: Compute the archetypes given \mathbf{B} , i.e. $\mathbf{Z} = \mathbf{B}\mathbf{X}$
 - 11: **end while**
 - 12: **return** $\mathbf{A}, \mathbf{B}, \mathbf{Z}$
-

The idea is to use a projected gradient algorithm to solve the objective in Equation (4).

First, we recast the optimization problem in terms of the l1-normalization invariant variables \tilde{a}_n and \tilde{b}_k (called invariant because these variables won't change if one applies l1-normalization)

$$\tilde{a}_{nk} = \frac{a_{nk}}{\sum_{k''=1}^K a_{nk''}}, \quad \tilde{b}_{kn} = \frac{b_{kn}}{\sum_{n''=1}^N b_{kn''}} \quad (25)$$

Then the gradient of the RSS wrt to a_n is obtained using the chain rule which yields

$$\begin{aligned} \frac{\partial \text{RSS}}{\partial a_n} &= \frac{\partial \text{RSS}}{\partial \tilde{a}_n} \frac{\partial \tilde{a}_n}{\partial a_n} \\ &= \left(\tilde{g}_n^{(A)} \right)^T \left(\frac{\left(\sum_{k''=1}^K a_{nk''} \right) \mathbf{I}_K - a_n \mathbf{1}_K^T}{\left(\sum_{k''=1}^K a_{nk''} \right)^2} \right) \\ &= \frac{\left(\sum_{k''=1}^K a_{nk''} \right) \left(\tilde{g}_n^{(A)} \right)^T \mathbf{I}_K - \left(\tilde{g}_n^{(A)} \right)^T a_n \mathbf{1}_K^T}{\left(\sum_{k''=1}^K a_{nk''} \right)^2} \end{aligned} \quad (26)$$

So for a single element we have

$$\begin{aligned}
\frac{\partial \text{RSS}}{\partial a_{nk}} &= \frac{\partial \text{RSS}}{\partial \tilde{a}_n} \frac{\partial \tilde{a}_n}{\partial a_{nk}} \\
&= \frac{\left(\sum_{k''=1}^K a_{nk''} \right) \tilde{g}_{nk}^{(A)} - \left(\tilde{g}_n^{(A)} \right)^T a_n}{\left(\sum_{k''=1}^K a_{nk''} \right)^2} \\
&= \frac{\left(\sum_{k''=1}^K a_{nk''} \right) \tilde{g}_{nk}^{(A)} - \sum_{k''=1}^K \tilde{g}_{nk''}^{(A)} a_{nk''}}{\left(\sum_{k''=1}^K a_{nk''} \right)^2}
\end{aligned} \tag{27}$$

If we additionally assume that a_n has been ℓ_1 normalized in the previous iteration we get

$$\frac{\partial \text{RSS}}{\partial a_{nk}} = \tilde{g}_{nk}^{(A)} - \sum_{k''=1}^K \tilde{g}_{nk''}^{(A)} a_{nk''} \tag{28}$$

which is exactly the same as in Section 2.2. of Mørup & Hansen (2012) [7]

To write down the algorithm we define P_{Σ_M} , a function that projects the rows of any matrix $\mathbf{H} \in \mathbb{R}^{N \times M}$ onto the M simplex

$$\begin{aligned}
\tilde{\mathbf{H}} &= P_{\Sigma_M}(\mathbf{H}) \quad \text{with} \\
\tilde{\mathbf{H}}_{nm} &= \frac{\max(\mathbf{H}_{nm}, 0)}{\sum_{m'=1}^M \max(\mathbf{H}_{nm'}, 0)}
\end{aligned} \tag{29}$$

Putting everything together, the algorithm in matrix notation is shown in Algorithm 2

3.4 Frank-Wolfe Algorithm

The idea of the Frank-Wolfe algorithm for archetypal analysis is to use gradient information, but to avoid the costly projection step of the PCHA.

As described above, the objective is convex in \mathbf{A} when fixing \mathbf{B} and vice versa. Furthermore, in this alternating optimization setting, the rows of \mathbf{A} and \mathbf{B} are constrained to the Σ_K and Σ_N simplex, respectively, which are convex sets. Thus, we have a convex minimization problem over a convex set which can be tackled using the efficient Frank-Wolfe algorithm [2]

4 Initialization

4.1 Furthest Sum

5 Weighted Archetypal Analysis

TODO

Algorithm 2 Principal Convex Hull Algorithm (PCHA)

```

1: Initialize  $\tilde{\mathbf{A}}, \tilde{\mathbf{B}}$ 
2: Initialize  $\mu_{\mathbf{A}} \leftarrow 1, \mu_{\mathbf{B}} \leftarrow 1$ 
3: while not converged or maximum number of iterations is reached do
4:   Update  $\mathbf{A}$  using projected gradient descent
5:    $\mathbf{Z} \leftarrow \tilde{\mathbf{B}}\mathbf{X}$ 
6:    $\text{RSS}_{\text{old}} \leftarrow \|\mathbf{X} - \mathbf{AZ}\|_F^2$ 
7:   for  $t = 1$  to  $T$  do
8:      $\tilde{\mathbf{G}}^{(\mathbf{A})} \leftarrow 2 \left( \tilde{\mathbf{A}}\mathbf{Z}\mathbf{Z}^T - \mathbf{X}\mathbf{Z}^T \right)$ 
9:      $\mathbf{G}^{(\mathbf{A})} \leftarrow \tilde{\mathbf{G}}^{(\mathbf{A})} - \left( \tilde{\mathbf{G}}^{(\mathbf{A})} \odot \mathbf{A} \right) \mathbf{1}_K \mathbf{1}_K^T$ 
10:    for  $j = 1$  to  $100T$  do ▷ line search
11:       $\mathbf{A} \leftarrow \mathbf{A} - \mu_{\mathbf{A}} \mathbf{G}^{(\mathbf{A})}$ 
12:       $\tilde{\mathbf{A}} \leftarrow P_{\Sigma_K}(\mathbf{A})$ 
13:       $\text{RSS}_{\text{new}} \leftarrow \|\mathbf{X} - \tilde{\mathbf{A}}\mathbf{Z}\|_F^2$ 
14:      if  $\text{RSS}_{\text{new}} < \text{RSS}_{\text{old}} + (1 + \epsilon)$  then
15:         $\mu_{\mathbf{A}} \leftarrow 1.2 \cdot \mu_{\mathbf{A}}$ 
16:        break
17:      else
18:         $\mu_{\mathbf{A}} \leftarrow 0.5 \cdot \mu_{\mathbf{A}}$ 
19:      end if
20:    end for
21:  end for
22:  Update  $\mathbf{B}$  using projected gradient descent
23:   $\text{RSS}_{\text{old}} \leftarrow \|\mathbf{X} - \mathbf{ABX}\|_F^2$ 
24:  for  $t = 1$  to  $T$  do
25:     $\tilde{\mathbf{G}}^{(\mathbf{B})} \leftarrow 2 \left( \tilde{\mathbf{A}}^T \tilde{\mathbf{A}} \tilde{\mathbf{B}}\mathbf{X}\mathbf{X}^T - \tilde{\mathbf{A}}^T \mathbf{X}\mathbf{X}^T \right)$ 
26:     $\mathbf{G}^{(\mathbf{B})} \leftarrow \tilde{\mathbf{G}}^{(\mathbf{B})} - \left( \tilde{\mathbf{G}}^{(\mathbf{B})} \odot \mathbf{B} \right) \mathbf{1}_N \mathbf{1}_N^T$ 
27:    for  $j = 1$  to  $100T$  do ▷ line search
28:       $\mathbf{B} \leftarrow \mathbf{B} - \mu_{\mathbf{B}} \mathbf{G}^{(\mathbf{B})}$ 
29:       $\tilde{\mathbf{B}} \leftarrow P_{\Sigma_N}(\mathbf{B})$ 
30:       $\text{RSS}_{\text{new}} \leftarrow \|\mathbf{X} - \tilde{\mathbf{A}}\tilde{\mathbf{B}}\mathbf{X}\|_F^2$ 
31:      if  $\text{RSS}_{\text{new}} < \text{RSS}_{\text{old}} + (1 + \epsilon)$  then
32:         $\mu_{\mathbf{B}} \leftarrow 1.2 \cdot \mu_{\mathbf{B}}$ 
33:        break
34:      else
35:         $\mu_{\mathbf{B}} \leftarrow 0.5 \cdot \mu_{\mathbf{B}}$ 
36:      end if
37:    end for
38:  end for
39:  Check for Convergence
40:   $\mathbf{Z} \leftarrow \tilde{\mathbf{B}}\mathbf{X}$ 
41:   $\text{RSS} \leftarrow \|\mathbf{X} - \tilde{\mathbf{A}}\mathbf{Z}\|_F^2$ 
42:  if RSS reduction is sufficient then
43:    break
44:  end if
45: end while
46: return  $\tilde{\mathbf{A}}, \tilde{\mathbf{B}}, \mathbf{Z}$ 

```

6 Dataset Size Reduction Methods

6.1 Coresets

TODO

7 References

- [1] Daniel Aloise et al. “NP-Hardness of Euclidean Sum-of-Squares Clustering”. In: *Machine Learning* 75.2 (May 2009), pp. 245–248. ISSN: 0885-6125, 1573-0565. DOI: 10.1007/s10994-009-5103-0. <http://link.springer.com/10.1007/s10994-009-5103-0> (visited on 12/07/2024). <http://link.springer.com/10.1007/s10994-009-5103-0>.
- [2] Kenneth L. Clarkson. “Coresets, Sparse Greedy Approximation, and the Frank-Wolfe Algorithm”. In: *ACM Trans. Algorithms* 6.4 (Sept. 3, 2010), 63:1–63:30. ISSN: 1549-6325. DOI: 10.1145/1824777.1824783. <https://doi.org/10.1145/1824777.1824783> (visited on 02/11/2025). <https://doi.org/10.1145/1824777.1824783>.
- [3] Adele Cutler and Leo Breiman. “Archetypal Analysis”. In: *Technometrics : a journal of statistics for the physical, chemical, and engineering sciences* 36.4 (1994), pp. 338–347. ISSN: 0040-1706. DOI: 10.1080/00401706.1994.10485840.
- [4] J. Eggert and E. Korner. “Sparse Coding and NMF”. In: *2004 IEEE International Joint Conference on Neural Networks (IEEE Cat. No.04CH37541)*. Vol. 4. July 2004, 2529–2533 vol.4. DOI: 10.1109/IJCNN.2004.1381036. <https://ieeexplore.ieee.org/document/1381036> (visited on 03/30/2025). <https://ieeexplore.ieee.org/document/1381036>.
- [5] Chih-Jen Lin. “Projected Gradient Methods for Nonnegative Matrix Factorization”. In: *Neural Computation* 19.10 (Oct. 2007), pp. 2756–2779. ISSN: 0899-7667. DOI: 10.1162/neco.2007.19.10.2756. <https://ieeexplore.ieee.org/document/6795860> (visited on 03/30/2025). <https://ieeexplore.ieee.org/document/6795860>.
- [6] Henryk Minc. *Nonnegative Matrices*. 2. print. Wiley-Interscience Series in Discrete Mathematics and Optimization. New York: Wiley, 1988. 206 pp. ISBN: 978-0-471-83966-8.
- [7] Morten Mørup and Lars Kai Hansen. “Archetypal Analysis for Machine Learning and Data Mining”. In: *Neurocomputing* 80 (Mar. 2012), pp. 54–63. ISSN: 09252312. DOI: 10.1016/j.neucom.2011.06.033. <https://linkinghub.elsevier.com/retrieve/pii/S0925231211006060> (visited on 12/07/2024). <https://linkinghub.elsevier.com/retrieve/pii/S0925231211006060>.
- [8] Kaare Brandt Petersen and Michael Syskind Pedersen. *The Matrix Cookbook*. Technical University of Denmark, 2012. <https://www.math.uwaterloo.ca/~hwolkowi/matrixcookbook.pdf>. <https://www.math.uwaterloo.ca/~hwolkowi/matrixcookbook.pdf>.

8 Appendix

8.1 Notation

- $N \in \mathbb{N}$ is the number of samples
- $D \in \mathbb{N}$ is the number of dimensions
- $K \leq \min(N, D)$ is the number of archetypes
- $\mathcal{X} = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}_{n=1}^N$ is our dataset, where each $\mathbf{x}_n \in \mathbb{R}^D$

- $\mathbf{X} \in \mathbb{R}^{N \times D}$ is our data matrix where each row is one sample
- $\mathbf{Z} \in \mathbb{R}^{K \times D}$ is our matrix of archetypes where each row is one archetype
- $\mathcal{Z} = \{\mathbf{z}_1, \dots, \mathbf{z}_K\}_{k=1}^K$ is set of archetypes, where each $\mathbf{z}_k \in \mathbb{R}^D$

8.2 Algorithms

Algorithm 3 Principal Convex Hull Algorithm (PCHA)

Require: Data matrix $\mathbf{X} \in \mathbb{R}^{N \times D}$, learning rates $\mu_{\mathbf{A}} > 0$, $\mu_{\mathbf{B}} > 0$

```

1: Initialize  $\mathbf{A}, \mathbf{B}$ 
2:  $\text{RSS}_{\text{old}} \leftarrow \|\mathbf{X} - \mathbf{A}\mathbf{B}\mathbf{X}\|_F^2$ 
3: while not converged do
  Update A coefficients:
4:    $\mathbf{Z} \leftarrow \mathbf{B}\mathbf{X}$  ▷ compute archetypes
5:    $\mathbf{G}^{(\mathbf{A})} \leftarrow 2(\mathbf{A}\mathbf{Z}\mathbf{Z}^T - \mathbf{X}\mathbf{Z}^T)$  ▷ gradient of RSS w.r.t.  $\mathbf{A}$ 
6:    $\mathbf{A} \leftarrow \mathbf{A} - \mu_{\mathbf{A}}\mathbf{G}^{(\mathbf{A})}$  ▷ gradient descent step
7:    $\mathbf{A} \leftarrow P_{\Sigma_K}(\mathbf{A})$  ▷ project rows of  $\mathbf{A}$  onto  $K$ -simplex
  Update B coefficients:
8:    $\mathbf{G}^{(\mathbf{B})} \leftarrow 2(\mathbf{A}^T\mathbf{A}\mathbf{B}\mathbf{X}\mathbf{X}^T - \mathbf{A}^T\mathbf{X}\mathbf{X}^T)$  ▷ gradient of RSS w.r.t.  $\mathbf{B}$ 
9:    $\mathbf{B} \leftarrow \mathbf{B} - \mu_{\mathbf{B}}\mathbf{G}^{(\mathbf{B})}$  ▷ gradient descent step
10:   $\mathbf{B} \leftarrow P_{\Sigma_N}(\mathbf{B})$  ▷ project rows of  $\mathbf{B}$  onto  $N$ -simplex
  Check convergence:
11:   $\text{RSS}_{\text{new}} \leftarrow \|\mathbf{X} - \mathbf{A}\mathbf{B}\mathbf{X}\|_F^2$ 
12:   $\text{rel\_decrease} \leftarrow \frac{\text{RSS}_{\text{old}} - \text{RSS}_{\text{new}}}{\text{RSS}_{\text{old}}}$  ▷ relative decrease in RSS
13:  if  $\text{rel\_decrease} < \epsilon$  then
14:    break ▷ convergence criterion met
15:  end if
16:   $\text{RSS}_{\text{old}} \leftarrow \text{RSS}_{\text{new}}$ 
17: end while
18: return  $\mathbf{A}, \mathbf{B}, \mathbf{Z}$ 

```
