

Convex non-negative matrix factorization for massive datasets

Christian Thureau · Kristian Kersting ·
Mirwaes Wahabzada · Christian Bauckhage

Received: 5 March 2010 / Revised: 28 July 2010 / Accepted: 9 October 2010 /
Published online: 26 October 2010
© Springer-Verlag London Limited 2010

Abstract Non-negative matrix factorization (NMF) has become a standard tool in data mining, information retrieval, and signal processing. It is used to factorize a non-negative data matrix into two non-negative matrix factors that contain basis elements and linear coefficients, respectively. Often, the columns of the first resulting factor are interpreted as “cluster centroids” of the input data, and the columns of the second factor are understood to contain cluster membership indicators. When analyzing data such as collections of gene expressions, documents, or images, it is often beneficial to ensure that the resulting cluster centroids are meaningful, for instance, by restricting them to be convex combinations of data points. However, known approaches to convex-NMF suffer from high computational costs and therefore hardly apply to large-scale data analysis problems. This paper presents a new framework for convex-NMF that allows for an efficient factorization of data matrices of millions of data points. Triggered by the simple observation that each data point can be expressed as a convex combination of vertices of the data convex hull, we require the basic factors to be vertices of the data convex hull. The benefits of *convex-hull NMF* are twofold. First, for a growing number of data points the expected size of the convex hull, i.e. the number of its vertices, grows much slower than the dataset. Second, distance preserving low-dimensional embeddings allow us to efficiently sample the convex hull and hence to quickly determine candidate vertices. Our extensive experimental evaluation on large datasets shows that convex-hull NMF compares favorably to convex-NMF in terms of both speed and reconstruction quality. We demonstrate that our method can easily be applied to large-scale, real-world datasets, in our case consisting of 750,000 DBLP entries, 4,000,000 digital images, and 150,000,000 votes on World of Warcraft® guilds, respectively.

Keywords Matrix factorization · Low-rank approximation · Data mining · Information retrieval · Large-scale data analysis

C. Thureau (✉) · K. Kersting · M. Wahabzada · C. Bauckhage
Fraunhofer IAIS, Schloss Birlinghoven, 53754 Sankt Augustin, Germany
e-mail: christian.thureau@iais.fraunhofer.de

1 Introduction

Matrix factorization is a fundamental step in many approaches to data mining, machine learning, and statistical pattern analysis. It can be found in application areas such as computational biology, computer vision, information retrieval, or social network analysis. Recent work in machine learning has focused on matrix factorizations which obey particular constraints that are inherent to certain data and therefore should be accounted for in any analysis. In particular, non-negative matrix factorization (NMF) focuses on the analysis of data matrices whose elements are non-negative, a common occurrence in representations of text or image data among others [4, 20]. Given a non-negative input matrix V , NMF aims at determining two non-negative matrix factors W and H such that

$$V \approx WH.$$

Convex non-negative matrix factorization (C-NMF) approaches additionally restrict the columns of W to be convex combinations of the data points gathered in V . This enforces W to represent meaningful “cluster centroids” and proves to be beneficial in applications such as text or genome mining, as well as image or social network analysis.

Although several (convex) NMF approaches have been proposed, very little work exists on how to apply NMF in the wild, i.e. in situations where one has to deal with millions of data points. This is exactly the problem we address in this paper.

Our main contribution is *convex-hull* non-negative matrix factorization (CH-NMF) which we demonstrate to be a very fast and scalable convex-NMF technique. Triggered by the simple observation that each data point is a convex combination of vertices of the data convex hull, the key idea is to even further restrict W and require it to only contain vertices of the convex hull. The benefits are twofold:

- First, the expected size of the convex hull typically grows considerably slower than that of the dataset. Consider for instance n random Gaussian points in the plane. Here, the expected number of vertices of the convex hull is $\Omega(\sqrt{\log n})$.
- Second, distance preserving, low-dimensional projections of the data allow one to efficiently sample the data convex hull. Candidate vertices can therefore be computed efficiently.

This article extends an earlier contribution [28]. In addition to related NMF algorithms, here we also compare the proposed method to other recent matrix factorization techniques such as the CUR or CMD decomposition introduced by [9] and [26]. Moreover, we discuss the merits of incorporating the FastMap algorithm for convex-hull sampling. Our extensive experimental evaluation on synthetic and large-scale real-world data shows that convex-hull NMF compares favorably to existing matrix factorization methods, in terms of both speed and reconstruction quality. Moreover, we show that convex-hull NMF can be easily applied to large-scale datasets. Specifically, we apply it to two large-scale, real-world datasets, consisting of more than 750,000 DBLP data records, of 4 million digital images, and of 150 million votes on World of Warcraft® guilds.

We proceed as follows: We begin by introducing definitions and notational conventions used in this article. Next, we will review NMF and related variants and then, in Sect. 4, we will develop convex-hull NMF. Section 5 will present our experimental evaluation on several synthetic and real-world datasets. Before concluding, we will briefly touch upon efficient methods for hull sampling.

2 Notation and definitions

Throughout this article, vectors are written as bold lower case letters (\mathbf{v}) and their entries are denoted using subscripted lower case italics (v_k). $\mathbf{0}$ is the vector of all zeros and $\mathbf{1}$ is the vector of all ones. We write $\mathbf{v} \geq \mathbf{0}$ to indicate that $v_k \geq 0$ for all k . The inner product of two vectors \mathbf{u} and \mathbf{v} is written as $\mathbf{u}^T \mathbf{v}$. Consequently, the expression $\mathbf{1}^T \mathbf{v}$ is a shorthand for $\sum_k v_k$.

Matrices are written using bold upper case letters (\mathbf{M}) and subscripted upper case italics (M_{ij}) denote individual entries. In order to indicate that \mathbf{M} is a real-valued $m \times n$ matrix, i.e. $\mathbf{M} \in \mathbb{R}^{m \times n}$, we may use the shorthand $\mathbf{M}^{m \times n}$. If the columns of a matrix are known, we also write $\mathbf{M} = [\mathbf{m}_1, \mathbf{m}_2, \dots, \mathbf{m}_n]$ where $\mathbf{m}_j \in \mathbb{R}^m$ is the j th column vector of \mathbf{M} .

In a slight abuse of notation, we may identify a finite set of vectors $\{\mathbf{v}_1, \dots, \mathbf{v}_n\}$ with the matrix $\mathbf{V} = [\mathbf{v}_1, \dots, \mathbf{v}_n]$ whose columns are given by the elements of the set. Moreover, we will simply use $\|\mathbf{M}\|$ to denote the *Frobenius norm* of \mathbf{M} so that $\|\mathbf{M}\|^2 = \sum_{i,j} M_{ij}^2$.

A set $S \subset \mathbb{R}^m$ is *convex*, if every point on the line segment between any two points in S is also in S , i.e. if $\forall \mathbf{u}, \mathbf{v} \in S, \forall \lambda \in [0, 1] : \lambda \mathbf{u} + (1 - \lambda) \mathbf{v} \in S$. A vector $\mathbf{v} \in \mathbb{R}^m$ is a *convex combination* of vectors $\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_l \in \mathbb{R}^m$, if $\mathbf{v} = \sum_i \lambda_i \mathbf{v}_i$ where $\lambda_i \geq 0$ and $\sum_i \lambda_i = 1$. Using matrix notation, we write convex combinations as $\mathbf{v} = \mathbf{V} \boldsymbol{\lambda}$ where $\mathbf{V} = [\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_l]$ and $\boldsymbol{\lambda} \in \mathbb{R}^l$ such that $\mathbf{1}^T \boldsymbol{\lambda} = 1$ and $\boldsymbol{\lambda} \geq \mathbf{0}$.

An *extreme point* of a convex set S is any point $\mathbf{v} \in S$ that is not a convex combination of other points in S . That is, if \mathbf{v} is an extreme point and $\mathbf{v} = \lambda \mathbf{u} + (1 - \lambda) \mathbf{w}$ for $\mathbf{u}, \mathbf{w} \in S$ and $\lambda \in [0, 1]$, then $\mathbf{v} = \mathbf{u} = \mathbf{w}$.

The *convex hull* \mathcal{C} of a set $S \subset \mathbb{R}^m$ is the set of all convex combinations of points in S , that is

$$\mathcal{C}(S) = \left\{ \sum_{\mathbf{v}_i \in \mathcal{R}} \lambda_i \mathbf{v}_i \mid \mathcal{R} \subseteq S, |\mathcal{R}| < \infty, \sum_i \lambda_i = 1, \lambda_i \geq 0 \forall i \right\}.$$

A *polytope* is the convex hull of finitely many points, i.e. it is the set $\mathcal{C}(S)$ for $|S| < \infty$. The extreme points of a polytope are called *vertices*. We use $\mathcal{V}(S)$ to denote the set of all vertices of a polytope. Note that every point inside a polytope can be expressed as a convex combination of the points in $\mathcal{V}(S)$.

3 Non-negative matrix factorization

Assume an $m \times n$ input data matrix $\mathbf{V} = [\mathbf{v}_1, \dots, \mathbf{v}_n] \in \mathbb{R}^{m \times n}$ that consists of n real-valued column vectors of dimensionality m . We consider factorizations of the form

$$\mathbf{V} \approx \mathbf{W} \mathbf{H}$$

where the matrix $\mathbf{W} \in \mathbb{R}^{m \times k}$ contains a set of $k \ll n$ basis vectors that are linearly combined using the coefficients in $\mathbf{H} \in \mathbb{R}^{k \times n}$ to represent the data.

Common approaches to achieve such a factorization include principal component analysis [15], singular value decomposition [11], vector quantization, or non-negative Matrix Factorization [23, 19]. Note that these methods impose different constraints and thus yield different matrix factors: principal component analysis constrains \mathbf{W} to be composed of orthonormal vectors and typically produces a dense \mathbf{H} , vector quantization constrains \mathbf{H} to unary vectors, and non-negative matrix factorization (NMF) assumes \mathbf{V} , \mathbf{W} , and \mathbf{H} to be non-negative matrices and often leads to sparse representations of the data. Therefore, in the case of NMF,

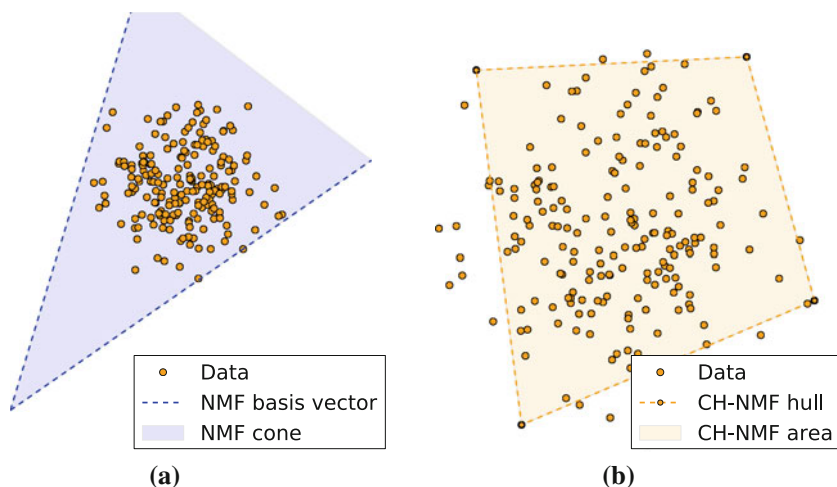


Fig. 1 Geometric interpretations of different variants of NMF. **a** The non-negativity constraints of standard NMF enforce the identification of the convex cone in which the data reside. **b** CH-NMF as proposed in this article factorizes the data such that the resulting basis vectors approximate the data convex hull

the factors in the matrix \mathbf{W} are often easy to interpret and \mathbf{H} tends to be sparse. From this point of view, NMF marks a middle ground between densely distributed and unary representations.

In addition to the data-compression aspects of NMF, the intuitive interpretability of the resulting factors makes it especially interesting for data mining. Figure 1a shows a geometric interpretation of NMF. Donoho and Stodden [8] pointed out that non-negative data reside in a convex cone and that a suitable set of non-negative basis vectors would span this cone. In fact, Meyer and Albright [18] investigated several NMF algorithms and found them to yield comparable solutions (up to scale) that typically represent points on the boundary of the cone. Hoyer [13] demonstrated that incorporating sparseness constraints leads to results that more closely resemble existing data points. We also note that Vasiloglou et al. [30] recently proved that for any choice of the number k of basis vectors, there exists a unique optimal solution; however, an algorithm that is guaranteed to find this solution in reasonable time is not yet available.

NMF has become an active area of research, and various variants and improvements have been proposed. For example, Cai et al. [3] present a matrix factorization that attempts to capture the geometric structure of the data. Kim and Park [16] accelerate NMF using projected gradient and alternating nonnegative least squares algorithms. Another interesting variation has recently been proposed by Suvrit [27] who bases the optimization steps in NMF on a block-iterative acceleration technique.

In this work, we focus on the problem of analyzing massive datasets using NMF. We build our approach on the idea of *Convex-NMF* (*C-NMF*) which was recently introduced by Ding et al. [7]. Convex-NMF aims at approximating the data matrix \mathbf{V} by means of convex combinations of data points, i.e.

$$\mathbf{V} = \mathbf{V}\mathbf{G}\mathbf{H}^T$$

where each column i of \mathbf{G} is a stochastic vector that obeys

$$\mathbf{1}^T \mathbf{g}_i = 1, \mathbf{g}_i \geq \mathbf{0}.$$

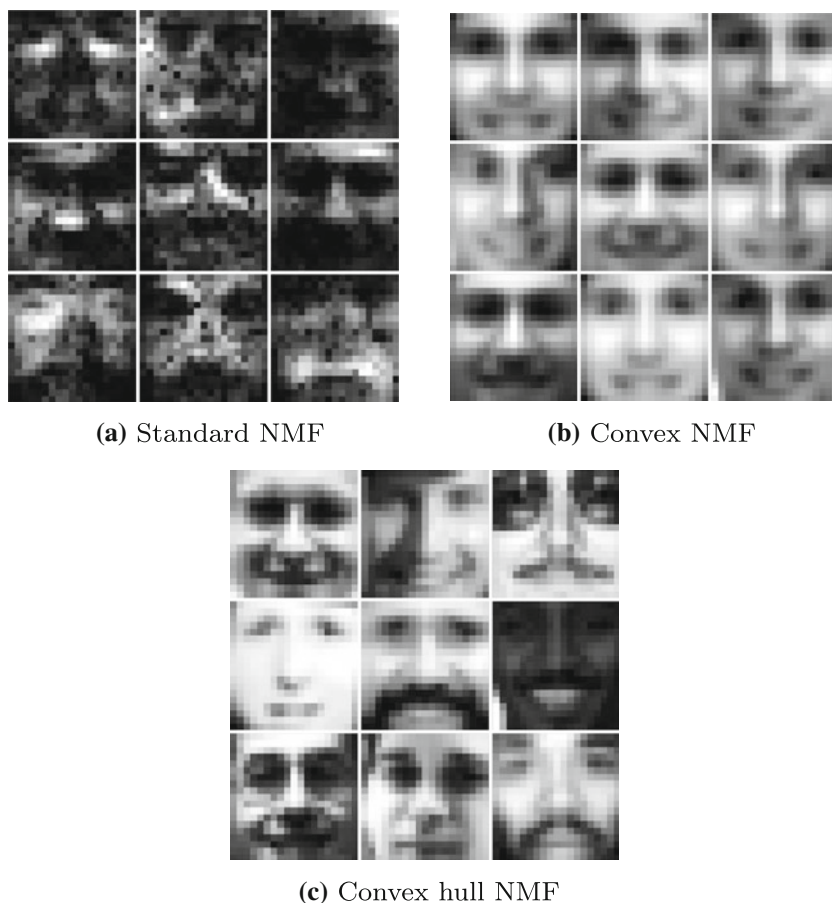


Fig. 2 Basis vectors resulting from different NMF variants applied to the CBCL Face Database. **a** Standard NMF results in sparse representations. Data points cannot be expressed as convex combinations of these basis elements. **b** C-NMF yields basis elements that allow for convex combinations. Moreover, the basis vectors are “meaningful” since they closely resemble given data points. They are, however, not indicative of characteristic variations among individual samples. **c** CH-NMF as proposed in this paper results in non-negative factors that represent such variations in the data (e.g. pale faces, faces with glasses, faces with beards, etc.)

This is akin to *Archetypal Analysis* according to Cutler and Breiman [5] where both matrices \mathbf{G} and \mathbf{H}^T are to be stochastic. Convex-NMF yields interesting interpretations of the data because each data point can now be expressed as a weighted sum of given data points. Consider the examples in Fig. 2 where we applied several NMF variants to analyze the CBCL Face Database 1 which consists of 2,429 19×19 gray-scale face images.¹ As one can see, standard NMF results in part-based, sparse representations. Data points, however, do not correspond to convex combinations of these elementary parts. C-NMF, on the other hand, yields basis elements that allow for expressing data points as convex combinations of given data points. Accordingly, the “meaning” of these basis elements is intuitively understandable.

¹ <http://cbcl.mit.edu/projects/cbcl/software-datasets/>.

In the remainder of this section, we will briefly review C-NMF. In particular, we shall point out its runtime complexity and relevancy for data analysis. C-NMF as introduced by Ding et al. [7] minimizes

$$J = \|V - VGH^T\|^2$$

where $V \in \mathbb{R}^{m \times n}$, $G \in \mathbb{R}^{n \times k}$, and $H \in \mathbb{R}^{n \times k}$. The matrices G and H are updated iteratively until convergence using the following update rules:

$$G_{ik} = G_{ik} \sqrt{\frac{(Y^+H)_{ik} + (Y^-GH^TH)_{ik}}{(Y^-H)_{ik} + (Y^+GH^TH)_{ik}}} \quad (1)$$

$$H_{ik} = H_{ik} \sqrt{\frac{(Y^+G)_{ik} + (HG^TY^-G)_{ik}}{(Y^-G)_{ik} + (HG^TY^+G)_{ik}}} \quad (2)$$

where $Y = V^TV$, and the matrices Y^+ and Y^- are given by

$$Y_{ik}^+ = \frac{1}{2}|Y_{ik}| + Y_{ik}$$

and

$$Y_{ik}^- = \frac{1}{2}|Y_{ik}| - Y_{ik},$$

respectively.

Two methods are proposed for the initialization of G and H . The first one initializes both matrices to (almost) unary representations based on a k-means clustering of V . The second one assumes a given NMF or Semi-NMF solution which is then refined using the above iterative algorithm. Note that in our experimental evaluation in Sect. 5 of this paper, we only present results using the first initialization scheme since we did not find the two alternatives to yield significantly different reconstruction errors.

Convex-NMF is closely related to k-means clustering and results in similar basis vectors W . However, it usually outperforms the k-means algorithm w.r.t. cluster accuracy.

Looking at the C-NMF update rules in (1) and (2), it becomes clear that they have a time complexity of $O(n^2)$ where n is the number of data point contained in V . Therefore, although there are only simple matrix multiplications at the core of the above algorithm, the size of the involved matrices quickly becomes a limiting factor, since V^TV results in an $n \times n$ matrix. Switching to an online update rule would avoid memory issues but it would at the same time introduce additional computational overhead. Overall, we can say that C-NMF does not scale to large datasets. In the following, we will present *Convex-Hull NMF (CH-NMF)*, which is a novel C-NMF method that is better suited for large-scale data analysis.

4 Convex-Hull NMF

Convex-Hull NMF aims at a factorization that incorporates data points that reside on the data convex hull. The corresponding data reconstruction has two interesting properties:

- First, the basis vectors correspond to given data points and mark, unlike in most other clustering or factorization techniques, the most extreme and not the most average data points.
- Second, any data point can be expressed as a convex and meaningful combination of these basis elements.

Both these characteristics offer interesting new opportunities for data interpretation. We exemplify this in Fig. 2 and will provide further evidence in Sect. 5.

Following Ding et al. [7], we consider a non-negative factorization of the form

$$\mathbf{V} = \mathbf{V}\mathbf{G}\mathbf{H}^T,$$

where $\mathbf{V} \in \mathbb{R}^{m \times n}$, $\mathbf{G} \in \mathbb{R}^{n \times k}$, $\mathbf{H} \in \mathbb{R}^{n \times k}$. Extending the approach by Ding et al. [7], we further restrict the columns of \mathbf{G} and \mathbf{H} to convexity, i.e.,

$$\begin{aligned} \mathbf{1}^T \mathbf{g}_i &= 1, & \mathbf{g}_i &\geq \mathbf{0} \\ \mathbf{1}^T \mathbf{h}_j &= 1, & \mathbf{h}_j &\geq \mathbf{0}. \end{aligned}$$

We note again that Ding et al. [7] also consider convex combinations but not for the matrix \mathbf{H} . In other words, in contrast to C-NMF, our idea of CH-NMF aims at factorizing the data such that each data point can be expressed as a convex combination of convex combinations of specific data points.

The task now is to solve the following constrained quadratic problem

$$\begin{aligned} \text{minimize } J &= \|\mathbf{V} - \mathbf{V}\mathbf{G}\mathbf{H}^T\|^2 \\ \text{subject to } \mathbf{1}^T \mathbf{g}_i &= 1, & \mathbf{g}_i &\geq \mathbf{0} \\ \mathbf{1}^T \mathbf{h}_j &= 1, & \mathbf{h}_j &\geq \mathbf{0}. \end{aligned} \quad (3)$$

In the following, we set

$$\mathbf{X} = \mathbf{V}\mathbf{G}$$

so that $\mathbf{X} \in \mathbb{R}^{m \times k}$. The intuition behind our algorithm is as follows: due to the above properties of the coefficient matrix \mathbf{G} , the column vectors in \mathbf{X} are convex combinations of columns in \mathbf{V} . Accordingly, by definition of the convex hull, the convex hull $\mathcal{C}(\mathbf{V})$ of \mathbf{V} must contain \mathbf{X} . Obviously, we could achieve a perfect factorization of the data matrix, giving $J = 0$ in (3), by choosing the columns of \mathbf{G} such that they would contain exactly one entry equal to 1 for each data point that is a vertex of the convex hull, while all other entries were set to zero. Or more informal: *following the definition of the convex hull, we can perfectly reconstruct any data point using a convex combination of data points that are vertices of the data convex hull.*

Therefore, our goal becomes to solve (3) by finding k appropriate data points on the convex hull of \mathbf{V} . In other words, we aim at solving

$$\begin{aligned} \text{minimize } J &= \|\mathbf{V} - \mathbf{X}\mathbf{H}^T\|^2 \\ \text{subject to } \mathbf{x}_i &\in \mathcal{V}(\mathbf{V}), \quad i = 1, \dots, k. \end{aligned} \quad (4)$$

Finding a solution to (4) is not necessarily straight forward. Rather, it is known that the worst-case complexity for computing the convex hull of n data points in m dimensions is $\Theta(n^{\frac{m}{2}})$. Consequently, computing the vertices of the convex hull of large datasets quickly becomes impractical. In this paper, we therefore propose an approximate solution that subsamples the convex hull but still offers convenient data reconstruction.

Our approach exploits the fact that any data point on the convex hull of a linear lower dimensional projection of the data also resides on the convex hull in the original data dimension. Formally, since \mathbf{V} contains finitely many points and therefore forms a polytope in \mathbb{R}^m , we can resort to

Algorithm 1 Convex-hull NMF using eigenvector projection as a mechanism for subsampling the convex hull.

- 1: Compute d eigenvectors $\mathbf{e}_1, \dots, \mathbf{e}_d$ of the covariance matrix derived from the data matrix $\mathbf{V}^{m \times n}$
 2: Project \mathbf{V} onto the 2D-subspaces

$$\mathbf{E}_{o,q}^{2 \times n} = \mathbf{V}^T [\mathbf{e}_o, \mathbf{e}_q], \quad o = 1 \dots d, \quad q = 1 \dots d, \quad o \neq q$$

- 3: Compute and mark convex hull vertices $\mathcal{C}(\mathbf{E}_{o,q})$ for each 2D projection
 4: Collect marked convex hull vertices (using the original data dimensionality m)

$$\mathbf{S}^{m \times p} = \mathcal{C}(\mathbf{E}_{1,2}) \cup \dots \cup \mathcal{C}(\mathbf{E}_{d-1,d})$$

- 5: Solve

$$\begin{aligned} \min \quad & J_S = \left\| \mathbf{S} - \mathbf{S} \mathbf{I}^{p \times k} \mathbf{J}^{k \times p} \right\|^2 \\ \text{s.t.} \quad & \mathbf{1}^T \mathbf{i}_i = 1, \quad \mathbf{i}_i \geq \mathbf{0} \\ & \mathbf{1}^T \mathbf{j}_i = 1, \quad \mathbf{j}_i \geq \mathbf{0} \end{aligned}$$

- 6: Compute

$$\mathbf{X}^{m \times k} = \mathbf{S} \mathbf{I}$$

- 7: Solve

$$\begin{aligned} \min \quad & J = \left\| \mathbf{v}_i - \mathbf{X} \mathbf{h}_i^T \right\|^2 \\ \text{s.t.} \quad & \mathbf{1}^T \mathbf{h}_i = 1, \quad \mathbf{h}_i \geq \mathbf{0} \end{aligned}$$

for $i = 1, \dots, n$

Theorem 1 (Main Theorem of Polytope Theory) *Every image of a polytope P under an affine map $\pi : \mathbf{x} \rightarrow \mathbf{M}\mathbf{x} + \mathbf{t}$ is a polytope.*

In particular, every vertex of an affine image of P , i.e., every vertex of the convex hull of the image of P , corresponds to a vertex of P . A proof of this important result can be found in the standard text by Ziegler [31].

Therefore, computing the vertices of the convex hulls of several 2D affine projections of the data offers a way of subsampling $\mathcal{V}(\mathbf{V})$. This is in fact an efficient way of doing so since computing the extreme points of the convex hull of a set of 2D points can be done in $O(n \log n)$ time [6].

This subsampling strategy is the main idea underlying Convex-Hull NMF, and its soundness directly follows from Theorem 1. Moreover, various methods can be used for linearly projecting the data to a 2D space:

- (a) **complete projections:** for lower dimensional data, it is possible to perform any pairwise projection of any data dimension.
- (b) **random projections:** randomly select the dimensions to be projected.
- (c) **eigenvector projections:** project the data using pairwise combinations of the first d eigenvectors of the covariance matrix of \mathbf{V} .
- (d) **FastMap projections:** use fastmap projections as introduced by Faloutsos and Lin [10].

Choosing PCA for linearly projecting the data to a 2D space, i.e. following option (c), leads to the Convex-Hull NMF approach that is summarized in Alg. 1. The mean and covariance matrix of \mathbf{V} can be computed iteratively, and the resulting matrices of size $m \times m$ and

can be efficiently stored. Projecting the data onto the $j = \frac{d(d-1)}{2}$ 2D subspaces resulting from pairwise combinations of the first d eigenvectors of the covariance matrix of \mathbf{V} copes with major variations in the data; in our experiments, we select d such that the first d corresponding eigenvalues account for 95% of the energy of the eigenvalue spectrum. Estimates of the resulting number of sampled points can be obtained from a result by Hueter [14], who showed that the expected size of the convex hull of n Gaussian data points in the plane is $\Omega(\sqrt{\log n})$. For Gaussian data, we thus expect to sample $p = j\sqrt{\log n}$ points; for data that can be approximated using a mixture of q Gaussian, we expect to sample $p = jq\sqrt{\log(n/q)}$ data points. In both cases, the candidate set grows much slower than n .

Given a candidate set $\mathbf{S} \subset \mathcal{V}(\mathbf{V})$ containing p vertices of the data convex hull, we now select those $k < p$ vertices that yield the best reconstruction of the remaining points in \mathbf{S} . Since \mathbf{S} is a $m \times p$ matrix, this, too, can be formulated as a NMF optimization problem with convexity constraints

$$\begin{aligned} &\text{minimize} && J_S = \|\mathbf{S} - \mathbf{S}\mathbf{I}\mathbf{J}\|^2 \\ &\text{subject to} && \mathbf{1}^T \mathbf{i}_i = 1, \quad \mathbf{i}_i \succeq \mathbf{0} \\ &&& \mathbf{1}^T \mathbf{j}_i = 1, \quad \mathbf{j}_i \succeq \mathbf{0} \end{aligned} \quad (5)$$

where $\mathbf{I} \in \mathbb{R}^{p \times k}$ and $\mathbf{J} \in \mathbb{R}^{k \times p}$. Since $p \ll n$, solving (5) can be done efficiently using common quadratic programming routines. Note that the data dimensionality is m . The convex hull projection only served to determine a candidate set; all further computations are carried out in the original data space.

Once a suitable $\mathbf{I} \in \mathbb{R}^{p \times k}$ has been determined, the matrix \mathbf{X} in (4) can be written as $\mathbf{X} = \mathbf{S}\mathbf{I}$ which guarantees that the minimization problem in (4) is solely concerned with k data points on the convex hull of \mathbf{V} . We found that \mathbf{I} usually results in unary representations. If this is not the case, we simply map $\mathbf{S}\mathbf{I}$ to their nearest neighboring data point in \mathbf{S} .

Given \mathbf{X} , the computation of the coefficients \mathbf{H} in (4) is straight forward. For smaller datasets, it is possible to use the iterative update rule from (2). However, since we do not further modify the basis vectors \mathbf{X} , we can also find an optimal solution for each data point \mathbf{v}_i individually by minimizing $J_i = \|\mathbf{v}_i - \mathbf{X}\mathbf{h}_i\|^2$, $\mathbf{1}^T \mathbf{h}_i = 1$, $\mathbf{h}_i \succeq \mathbf{0}$ using common solvers. Obviously, this can be easily parallelized.

5 Experimental evaluation

In this section, we present experimental evaluations of our CH-NMF algorithm on different datasets. Our first experimental setup compares CH-NMF to related methods. We evaluate the runtime behavior and reconstruction accuracy of CH-NMF and C-NMF, and we compare the proposed projection sampling to common random sampling techniques. For measuring the reconstruction accuracy, we usually compute the Frobenius norm or the relative Frobenius norm $R = \|\mathbf{V} - \mathbf{V}'\|^2 / \|\mathbf{V}\|^2$, where \mathbf{V}' denotes the reconstruction of the data matrix \mathbf{V} using the estimated basis vectors and coefficients (the accuracy is then defined as: accuracy = $1 - R$). We consider synthetic data and several small real-world datasets. Evaluation of C-NMF and other factorization methods is rather infeasible for very large datasets because these methods scale with $O(n^2)$, where n is the number of data. Thus, we limit the maximum number of randomly generated data to 5,000 points and use real-world datasets of even smaller size.

Since our main interest is in the analysis of very large real-world datasets, we also present experiments that evaluate CH-NMF on three challenging datasets.

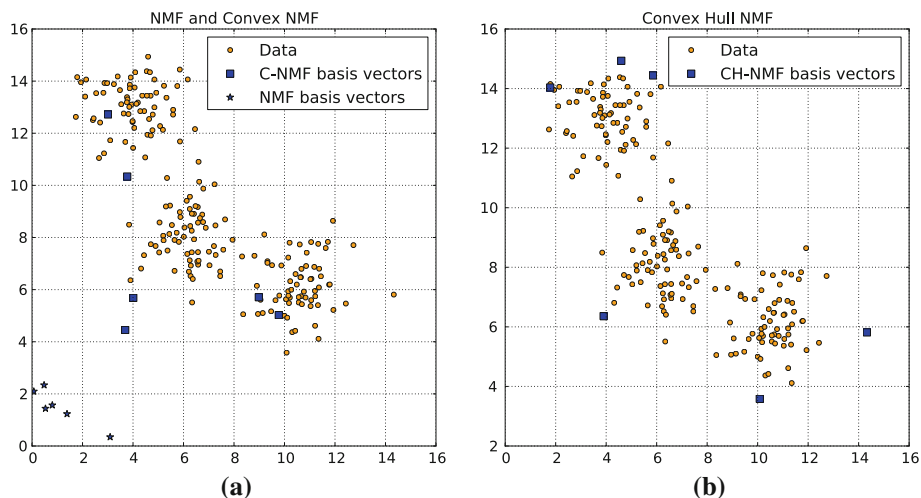


Fig. 3 Basis vectors produced by convex-NMF, convex-hull-NMF, and NMF for data sampled from randomly placed Gaussian distributions in 2D. It can be seen that CH-NMF basis vectors reside on the convex hull of the data, whereas C-NMF usually converges against cluster centroids. NMF basis vectors, in contrast, do not correspond to actual data points; however, they could be rescaled to reside closer to actual data points

- The first real-world dataset is a snapshot of the DBLP citation database² and we consider 757,368 data vectors in \mathbb{R}^{68} .
- Our second large dataset contains 70-dimensional feature vectors and comprises 150,000,000 player activity scores for the Massively Multiplayer Online Game World of Warcraft[®].
- The third massive dataset consists of 4,000,000 images of the Tiny Image dataset collected by [29]. The images are RGB color pictures of size 32×32 which are represented by means of 384-dimensional GIST feature vectors [21].

Our experimental evaluation is based on a basic Python implementation³ that is running on a standard Intel-Quadcore 2. GHz computer. For the quadratic optimization steps in our implementation, we apply the `cvxopt` library provided by Dahl and Vandenberghe.⁴ Although CH-NMF can be easily parallelized, we used a strictly serial implementation.

5.1 Method Comparison

Similar to the experimental procedure by [7], we evaluated runtime performance using varying numbers of data points sampled from three randomly positioned 2D Gaussians (see Fig. 3). In addition to measuring computation times, we also computed the mean reconstruction accuracy.

First of all, we compared CH-NMF against C-NMF. As an additional baseline, and since it is closely connected to C-NMF, we also evaluated k-means clustering. We varied the number of data points ranging from 250 to 5,500 in steps of 250. In order to account for effects of

² <http://www.informatik.uni-trier.de/~ley/db/>.

³ <http://www-kd.iai.uni-bonn.de/index.php?page=software>.

⁴ <http://abel.ee.ucla.edu/cvxopt/>.

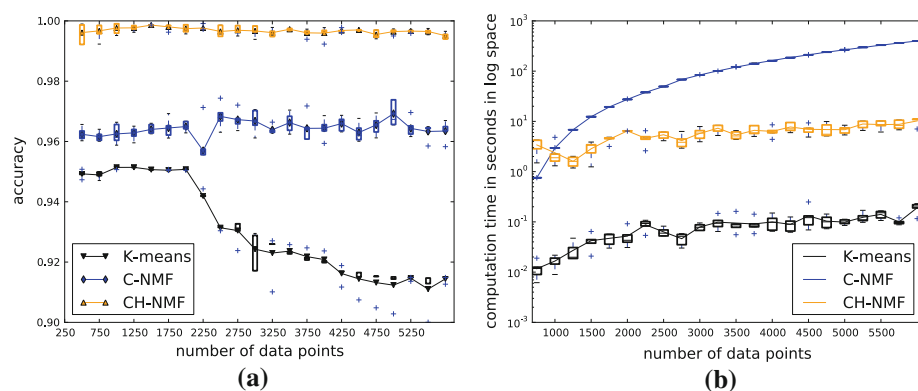


Fig. 4 **a** Boxplots of reconstruction accuracy of k-means, C-NMF, and CH-NMF for varying numbers of synthetically generated data. As one can see, CH-NMF outperforms both other methods. **b** Boxplots of computation time (in log-space) for k-means, C-NMF, and CH-NMF for varying numbers of synthetically generated data. It can be seen that, for smaller datasets of up to 1,000 samples, CH-NMF is actually slower than C-NMF. However, for sample size larger than 1,000 CH-NMF shows significant speed-ups. The computation time remains almost constant as it mainly grows with the number of vertices of the convex hull and not with the size of the data

random initializations, we repeated each experiment 5 times. The number of basis vectors (or clusters for k-means) was set to 6. The resulting average reconstruction accuracies can be seen in Fig. 4a, and the observed runtime performance is shown in Fig. 4b. Following Ding et al. [7], we carried out 100 iterations for each approach.

Overall, CH-NMF provides a very accurate data reconstruction. This confirms the intuition that forms the basis of our work: using vertices of the convex hull in order to reconstruct data points must lead to a near perfect reconstruction. Since we exploit this property and subsample the convex hull, we obtain very accurate data reconstruction on average. In comparison, C-NMF often converges to basis vectors that correspond to local modes of the data. While this is by itself a desirable property, it yields a slightly worse reconstruction. We attribute this to the fact that C-NMF does impose a convexity constraint on the coefficients in \mathbf{H}^T . This results in a conic reconstruction that limits perfect reconstruction to data points within the conic hull of basis vectors [17].

Regarding runtime performance, it can be seen that while CH-NMF initially takes longer for fewer data samples (up to 1,000), its computation time increases only moderately for a growing number of data. This is due to the moderate growth of the number of vertices of the convex hull w.r.t. the total number of data points [14]. For n data points that are distributed according to a mixture of q Gaussians, we usually have far less than $\Omega(q\sqrt{\log n})$ data point on the convex hull. Accordingly, it only takes about 54 s to compute basis vectors from 50,000 data points. It should be noted that about 90% of this time is spent on the constrained quadratic minimization problems $J = \|\mathbf{v}_i - \mathbf{X}\mathbf{h}_i^T\|^2$ which could be solved in parallel. Memory requirements for CH-NMF are rather low as it only requires to store k convex-hull data points at maximum.

For comparing C-NMF and CH-NMF on real-world datasets, we use the CBCL Face Database 1 (2,429 19×19 gray-scale face images), and a random selection of 2,000 images from the already mentioned Tiny Image dataset [29]. The results are summarized in Table 1. It can be seen that overall CH-NMF and C-NMF yield a very similar performance for smaller datasets. Usually, CH-NMF is faster and has a better reconstruction accuracy.

Table 1 Comparison of C-NMF and CH-NMF for two smaller real-word datasets

# Basis vectors	5	10	20	30
<i>(a) CBCL Face database</i>				
Accuracy CH-NMF	0.783	0.802	0.82	0.829
Time in s CH-NMF	27.788	43.011	69.473	95.125
Accuracy C-NMF	0.771	0.78	0.786	0.788
Time in sec. C-NMF	48.121	59.906	90.131	120.756
<i>(b) 2000 random images from the Tiny Image dataset</i>				
Accuracy CH-NMF	0.564	0.595	0.616	0.626
Time in sec. CH-NMF	19.532	22.034	38.477	44.73
Accuracy C-NMF	0.563	0.578	0.586	0.593
Time in s C-NMF	35.996	43.832	64.084	84.102

As C-NMF can not be run on very large datasets, we compare the average accuracy and runtime on the CBCL Face dataset and a selection from the Tiny Image dataset. It can be seen that both approaches yield rather similar accuracy for varying numbers of basis vectors. As expected for small datasets, CH-NMF is only slightly faster

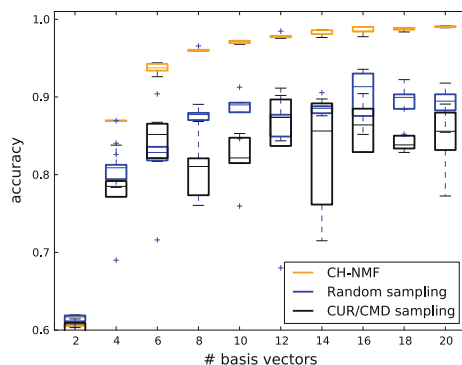
Table 2 Basis vector coefficients for a number of selected World of Warcraft[®] guilds (the corresponding basis vectors are explained in detail in Fig. 8)

	Basis vector coefficients							
	Fig. 8a	Fig. 8b	Fig. 8c	Fig. 8d	Fig. 8e	Fig. 8f	Fig. 8g	Fig. 8h
The legacy	0.313	0.185	0.33	0.119	0.013	0.0	0.04	0.0
Closure	0.0	0.4	0.325	0.0	0.0	0.275	0.0	0.0
Refuge	0.139	0.373	0.455	0.0	0.004	0.014	0.0	0.016
Risen	0.042	0.144	0.454	0.001	0.0	0.359	0.0	0.0
Wraith	0.306	0.0	0.31	0.23	0.089	0.0	0.065	0.0
Security	0.625	0.021	0.146	0.08	0.051	0.0	0.077	0.0
Ensidia	0.0	0.0	0.552	0.0	0.0	0.448	0.0	0.0
Apex	0.084	0.0	0.634	0.072	0.21	0.0	0.0	0.0
Premonition	0.16	0.013	0.545	0.0	0.181	0.0	0.102	0.0
For the Horde	0.053	0.0	0.396	0.2	0.35	0.0	0.0	0.0
Loot FTW	0.091	0.0	0.428	0.169	0.275	0.0	0.038	0.0
Method	0.044	0.0	0.485	0.412	0.059	0.0	0.0	0.0
Experience	0.036	0.318	0.555	0.0	0.0	0.058	0.007	0.026
Last resortk	0.119	0.499	0.202	0.007	0.0	0.174	0.0	0.0
Incorporated	0.0	0.237	0.417	0.0	0.0	0.346	0.0	0.0

These guilds are a selection from the top 20 guilds worldwide. Their coefficients in a basis representation resulting from CH-NMF tend to be similar; they all show a strong tendency toward the basis vector in Fig. 8c “formed after first update, very active with second update” just as it is indicated in the projection in Fig. 7

Apart from directly sampling the convex hull, one could try to reduce the computational complexity of very large datasets by employing random sampling strategies. The most simple approach would be a uniform random sampling, i.e. each data sample is picked randomly assuming a uniform distribution. While this leads to very fast sampling, the resulting

Fig. 5 Comparison of random sampling strategies to CH-NMF for the convexity constrained matrix factorization problem (averaged over 5 runs). While random sampling and CUR/CMD sometimes yield convenient results, CH-NMF is more stable and provides a higher accuracy



factorizations will be suboptimal as it effectively ignores large parts of the data. CUR matrix decomposition, as introduced by Drineas et al. [9], provides a smarter subsampling. In CUR decompositions, the columns (or rows) of a matrix are sampled from a weighted distribution. The sampling probability for picking a column (or row) depends on its vector norm, and columns with a higher norm (a high statistical leverage) have a higher sampling probability. *Compact Matrix Decomposition (CMD)* as proposed by Sun et al. [26] further improves on CUR by avoiding double selections and is therefore faster to compute. As CMD and CUR usually yield the same accuracy Sun et al. [26], we use them synonymously in the following.

For comparing the proposed convex-hull sampling to random sampling strategies, we changed the sampling to a uniform random sampling and CMD/CUR and then applied C-NMF to the reduced dataset. The resulting basis vectors are then used for computing coefficients for the complete dataset and measuring the reconstruction quality. We sampled 500 three-dimensional data points from 5 randomly placed Gaussian distributions. We first run CH-NMF using the eigenvector-based sampling. As CH-NMF automatically finds a number of suitable data points to sample, we could use this for specifying the number of random samples. For the relatively small sample sizes considered here, runtime performance for the three approaches is very similar. It should be noted that for very large datasets, uniform random sampling or CMD/CUR are faster than the CH-NMF sampling strategies.

Figure 5 shows the results for different numbers of basis vectors. Overall, it can be seen that CH-NMF provides the best reconstruction accuracy and is also more stable than the random sampling strategies. Interestingly, CMD/CUR show a worse performance than random sampling. The statistical leverage does not seem to be a suitable criterion for convexity constrained matrix factorization as it ignores possibly good candidates that reside on the convex hull but have a low statistical leverage.

Next, we will present experiments on large, real-world datasets that are impractical for established NMF methods but can be processed in reasonable time using the proposed CH-NMF technique.

5.2 Bibliographic analysis: DBLP

Bibliographic databases such as DBLP⁵ are a rich source of information. In the experiments reported here, we are interested in the question whether there are common patterns in the development of academic careers. To this aim, we extracted from DBLP the cumulative publication histograms of 757, 368 authors (see Fig. 6). A publication histogram consists of the

⁵ <http://www.informatik.uni-trier.de/~ley/db/>.

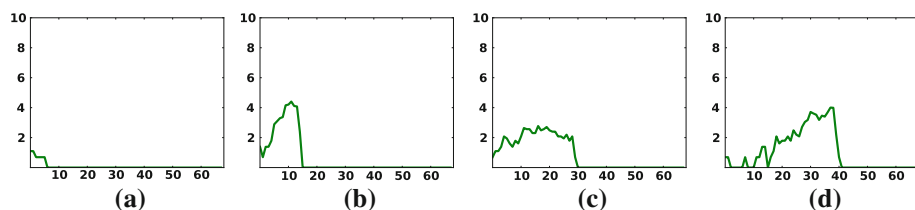


Fig. 6 Basis vectors resulting from the application of CH-NMF to publication histograms extracted from DBLP. Going from (a) to (d), the histograms reflect publication frequencies at different stages of an academic career. For example, **a** could describe a Ph.D. candidate that just started with a very few publications, and **d** describes a group leader who continuously increased his/her publication output

number of publications listed in DBLP in her first year, second year, and so on. We have cumulated the publications numbers of the years. The longest histogram we found spanned a period of 68 years. To get curves of equal length, we filled missing years with zeros. Following Aitchison [2], we use logarithmic histogram values in our analysis. The idea is that the publication histogram of an author provides a good characterization of her activity and also, to some, extend of her success (but of course does not suggest impact or quality).

Application of CH-NMF to publication histograms yields interesting basis histograms, which can be seen in Fig. 6. Again, it is important to note that the resulting basis vectors correspond to real data or, in this particular case, to real people. Other publication records can now be easily expressed as a convex combination of the most extreme academic career profiles.

5.3 Activity mining: World of Warcraft[®]

This dataset consists of recordings of the online appearance of characters in the computer game World of Warcraft[®]. To the best of our knowledge, this is the first time that vast recordings of in-game data of Massively Multiplayer Online games are considered as a source for data mining, we will briefly describe the data in more detail.

At the time of this writing, World of Warcraft[®] has more than 12 million paying customers. The game takes place in a virtual medieval fantasy environment. However, instead of one persistent world, players are distributed among different so-called realms. These realms exist in parallel and can have slightly varying rule-sets, i.e. each realm is its own world. In the United States and Europe, there exist about 500 unique realms. World of Warcraft[®] is often considered one large social platform which is used for chatting, team-play, and gathering. Compared to well-known virtual worlds that mainly serve as chat platforms such as, for example, Second-life[®], World of Warcraft[®] is probably the *real* second life as it has a larger and more active (paying) user base. Moreover, a whole industry is developing around World of Warcraft[®]. It is estimated that 400,000 people worldwide are employed as gold-farmers who are collecting virtual goods for online games and sell them over the Internet.

Players organize in groups, which are called guilds. Unlike groups known from other social platforms or services such as Facebook or Flickr, membership in a guild is exclusive. Obviously, their choice of a guild influences with whom players frequently interact. It also influences how successful players are in terms of in-game achievements, for instance, how likely they are to obtain better equipment or rare items.

Data we consider here were crawled from a publicly accessible Web site.⁶ In the following, we interpret each online appearance of a character as a vote. Character observations span

⁶ <http://www.warcraftrealms.com>.

Convex Non-negative Matrix Factorization for Massive Datasets

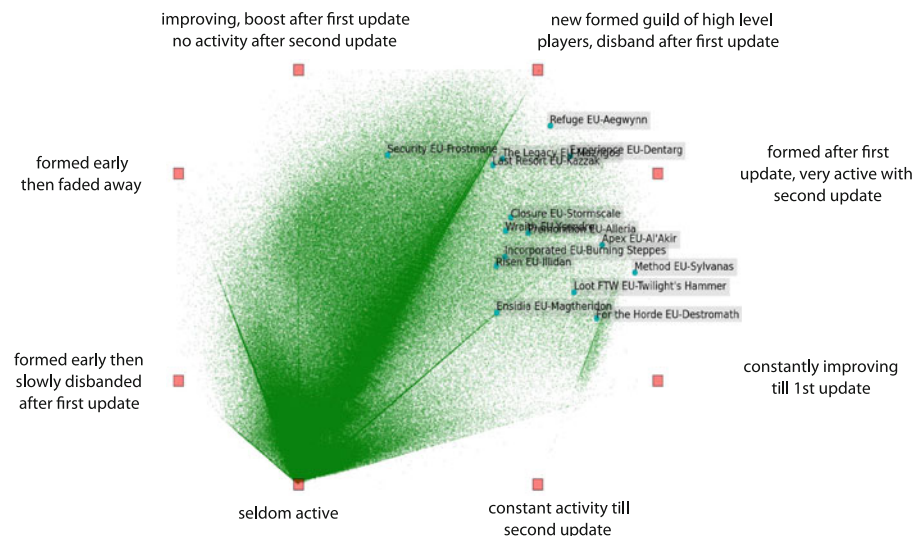


Fig. 7 The world of World of Warcraft[®]. The figure shows a 2D projection of 150 million votes for guilds into the space spanned by the convex-hull NMF bases (*boxes*). Intuitively understandable interpretations of the basis elements are shown right next to them. The names of 13 of the 20 most successful guilds (according to <http://www.wowprogress.com/>.) are shown inside the plot. Computing the factorization and the embedding took less than 1.5 h on a standard desktop machine

a period of 4 years. Every time, a character is seen online, (s)he votes for the guild (s)he is a member of; votes are cast according to a player's level. We accumulate the votes into a level-guild histogram, ranging from level 10 (levels 1–9 are excluded) to level 80 (the highest possible level). Players advance in level by engaging in in-game activity, i.e. by completing quests or other heroic deeds.

We assume that the level distribution among a guild is a good descriptor for its success. For example, a guild of very experienced level 80 characters has a higher chance of great achievements than a guild of level 10 players. Also, a level histogram provides an indicator of player activity over time. If players are continuously staying with a particular guild, we expect an equally distributed level histogram, as the characters increase their levels over time. Again following Aitchison [2], we use logarithmic histogram values in our analysis. In total, we collected 150 million votes of 18 million characters belonging to 1.4 million guilds.

Application of CH-NMF to the World of Warcraft[®] data reveals interesting structures in the data. Figure 7 shows a projection of all level-guild histograms into the space spanned by the CH-NMF basis vectors. We decided for 8 basis vectors since we assumed that one basis vector would function as a unity vector and encode a range of 10 levels. However, the actual outcome differs from what we expected. Figure 8 shows intuitive explanations of the basis vectors we obtained. They provide an interpretation of the distribution of all World of Warcraft[®] guilds; in particular, it is interesting to note that the majority of guilds are very close to the basis vector in Fig. 8f. Apparently, the vast majority of guilds or players seem to be casual gamers. Contrarily, none of the most successful guilds appear to belong to that group. Also, notice that we can spot particular events in the history of the game just from looking at the basis level-guild histograms. Considerable updates of the game content

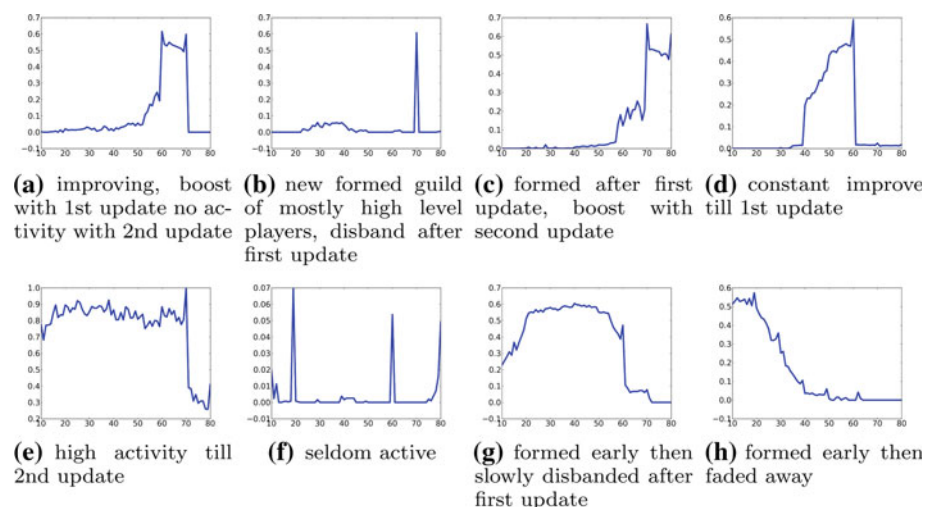


Fig. 8 Basis vectors resulting from applying CH-NMF to the World of Warcraft[®] data. Each subplot shows a level-guild histogram. The x -axis denotes histogram bins corresponding to player levels. Player levels describe a character's experience and may increase over time. Each bin contains accumulated player votes in log-space. By definition of CH-NMF, the basis vectors show the most extreme data points. Accordingly, one can easily provide intuitive explanations of the determined basis vectors. For example, **c** shows a guild that mainly consists of high level players whose experience ranges from level 70 to 80. Therefore, we can assume that the guild was formed after the 2nd update, when the level cap was raised to 80 and then stayed together. **h**, as a second example, shows high activity of players of lower levels (up to level 30). Here, we can assume that the guild was formed early after the release of the game and then disbanded

(a regular procedure that makes novel content available and also allows a further advancement in character level) appear to entail a restructuring of social groups.

Regarding runtime performance, the application of CH-NMF took 57 min (this includes the complete algorithm described in Fig. 1). The computation of 8 basis vectors only took 7.5 min. As expected, most time was spent on computing the 1.4 million basis vector coefficients and on the computation of level-guild histograms (about 20 min).

5.4 Image retrieval: large-scale image collection

The third massive dataset that we considered in our experiments consists of 4 million images downloaded from the Internet [29]. The images are rescaled to 32×32 pixels in 3 color channels and are represented by means of 384-dimensional GIST feature vectors. A selection of 25 basis vectors obtained from CH-NMF can be seen in Fig. 9. Interestingly, some of the basis vectors found among the tiny images bear a geometric similarity to 2D Fourier basis functions or Gabor filters. In fact, this stood to be expected, because the GIST features introduced by Oliva and Torralba [21] can be seen as a frequency domain representation of the input images. Accordingly, those images whose geometric structure is most similar to superpositions of elementary 2D sine or cosine functions form the extreme points in this space. Using larger numbers of basis vectors added more and more structured images to the set of basis images.

Due to the higher dimension of the data considered in this experiment, the application of CH-NMF took longer than in the case of the World of Warcraft[®] data. Using the idea of FastMap projections for an efficient pre-selection of candidate vectors for CH-NMF, it took

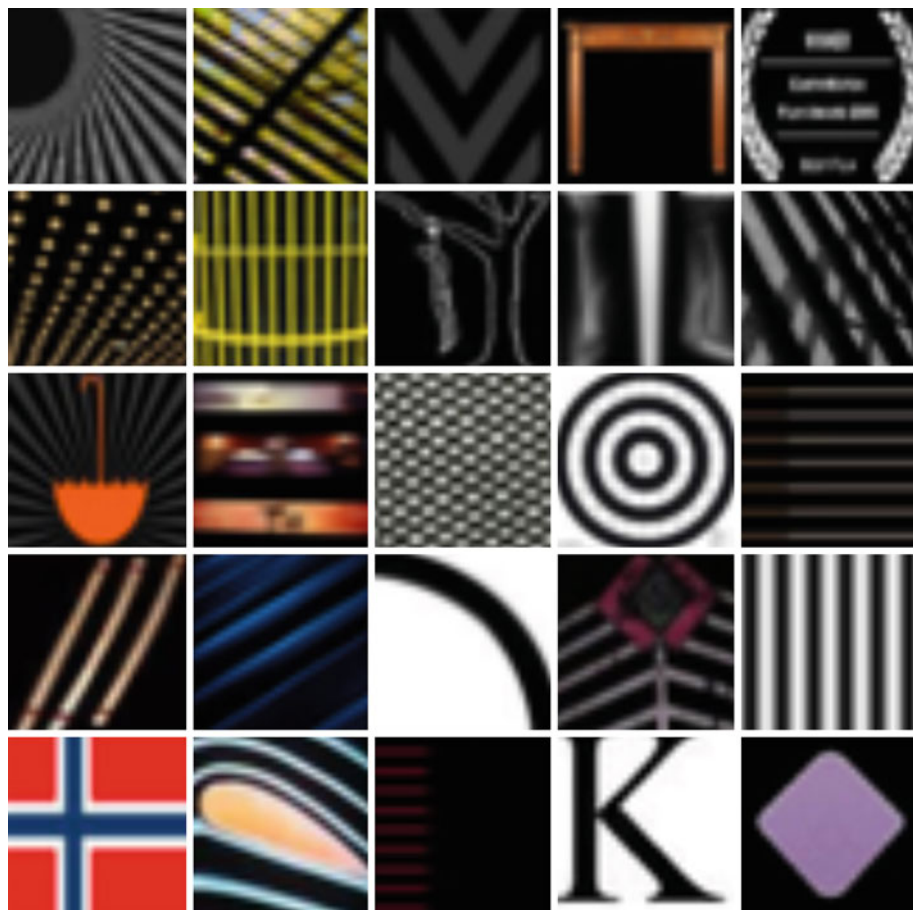


Fig. 9 Twenty-five basis vectors (shown as images) that result from applying CH-NMF to a dataset of 4,000,000 tiny images. Interestingly, the structure of these basis images resembles elementary Fourier basis functions or Gabor filters

about 44 h on a single computer to determine the 25 basis elements shown in Fig. 9. This is, in fact, a remarkable feat as the data were stored on an external USB hard drive, and runtime could be considerably reduced using an internal hard disk.

6 Sampling the convex hull

For massive, high-dimensional data, computing the covariance matrix can take a lot of time. In these cases, we use the FastMap projections to sample vertices from the data convex hull. We briefly summarize the algorithm in the following. FastMap computes a u -dimensional Euclidean embedding and proceeds as follows. Given pairwise distances among the columns of the data matrix V , we select a pair of distant vectors called *pivot objects*. The line segment between the pivot objects serves as the first coordinate axis. For each data point v , we determine its coordinate value along this axis by projecting v onto this line. Next, the



(a) Fastmap selection, 3 basis vectors



(b) PCA selection, 3 basis vectors



(c) Fastmap selection, 6 basis vectors



(d) PCA selection, 6 basis vectors



(e) Fastmap selection, 9 basis vector



(f) PCA selection, 9 basis vectors

Fig. 10 Comparison of basis vectors resulting from using Fastmap and PCA for sampling convex-hull vertices. While the resulting basis images differ between the two methods, for each image in the *left column* of the figure we can identify a very similar counterpart in the *right column*

pairwise distances of all objects are updated to reflect this projection, i.e. we compute the pairwise distances among the data points in the subspace orthogonal to the line. This process is repeated until, after u iterations, we obtain the u coordinates as well as the u -dimensional representation of all data points. Ostrouchov and Samatova [22] have shown that the pivots are taken from the faces, usually vertices, of the convex hull of the data points in the original implicit Euclidean space. This justifies the idea to employ FastMap in step (3) of the CH-NMF algorithm discussed in Sect. 4. An experimental comparison of the PCA and FastMap-based subsampling strategies is shown in Figs. 10 and 11.

First, we compared CH-NMF with PCA against CH-NMF with FastMap on synthetic data. We evaluated the mean reconstruction error and runtime performance (see Fig. 11) using ten

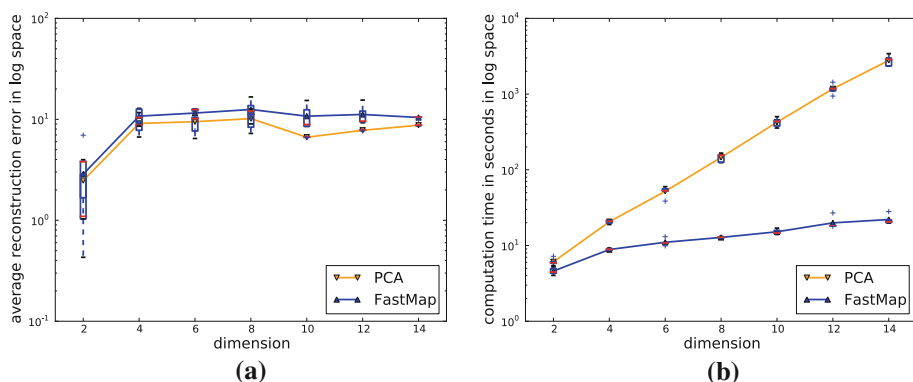


Fig. 11 Results of using different sampling strategies on synthetic data. While the average reconstruction error is almost identical for Fastmap and PCA (a), Fastmap subselection is generally much faster to compute (b)

randomly positioned Gaussians in 2D with varying number of dimensions (ranging between 2 and 14) and a total number of 4000 data points. We repeated each experiment 5 times to account for effects of random initialization. The number of basis vectors was set to the number of dimensions + 1. As can be seen, using Fastmap for candidate selection is not only several orders of magnitude faster than sampling based on PCA, it also manages to achieve comparable reconstruction errors. As a second test, we provide a qualitative evaluation of the resulting basis vectors. Using the CBCL face database, we considered candidate sets obtained from Fastmap and PCA for computing CH-NMF. Figure 10 shows the resulting bases. It can be seen that, while the resulting bases are not identical, they show a high degree of similarity.

7 Conclusions

Machine learning and data mining techniques typically consists of two parts: the model and the data. Most effort in recent years has gone into the modeling part. Large-scale datasets, however, allow one to move into the opposite direction [1, 12, 29] and ask the question of how much can the data itself help us to solve the problem. This direction is particularly appealing, given that the Internet nowadays offers a plentiful supply of large-scale datasets for many challenging tasks.

Motivated by this, we have presented a data-driven NMF approach, called *convex-hull NMF*, that is fast and scales well. The key idea is to restrict the “clusters” to be combinations of vertices of the convex hull of the dataset and thus directly explore the data itself in order to solve the convex-NMF problem. Our experimental results reveal that convex-hull NMF can indeed effectively extract meaningful “clusters” from datasets containing millions of images and rating.

For future work, it is interesting to apply convex-hull NMF to other challenging datasets, such as Wikipedia, Netflix, Facebook, or the blogosphere and to use it for applications such as collaborative filtering. To improve performance, one could explore mixtures of convex hulls to deal with general distributions. Also, an analysis as to how outliers affect the performance on small- and medium-sized datasets should be done, as our results indicate that this is not an issue for very large datasets. Finally, a comparison to max-margin matrix

factorization [25,24], which computes a low-norm instead of a low-rank factorization, would be interesting. Overall our experimental results are an encouraging sign that applying NMF techniques in the wild, i.e. on millions of data points may not be insurmountable.

Acknowledgments We thank the anonymous reviewers for their useful comments and suggestions. The authors would also like to thank A. Torralba, R. Fergus, and W. T. Freeman for making the tiny images freely available. Kristian Kersting and Mirwaes Wahabzada were supported by the Fraunhofer ATTRACT Fellowship STREAM.

References

1. Aggarwal C (2009) On classification and segmentation of massive audio data streams. *Knowl Inf Syst* 20(2):137–156
2. Aitchison J (1982) The statistical analysis of compositional data. *J R Stat Soc B* 44(2):139–177
3. Cai D, He X, Wu X, Han J (2008) Non-negative matrix factorization on manifold. In: *Proceedings of IEEE international conference on data mining*
4. Chen Y, Rege M, Dong M, Hua J (2008) Non-negative matrix factorization for semi-supervised data clustering. *Knowl Inf Syst* 17(3):355–379
5. Cutler A, Breiman L (1994) Archetypal analysis. *Technometrics* 36(4):338–347
6. de Berg M, van Kreveld M, Overmars M, Schwarzkopf O (2000) *Computational geometry*. Springer, Heidelberg
7. Ding C, Li T, Jordan M (2009) Convex and semi-nonnegative matrix factorizations. *IEEE Trans Pattern Anal Mach Intell* 32(1):45–55
8. Donoho D, Stodden V (2004) When does non-negative matrix factorization give a correct decomposition into parts?. In: *Advances in neural information processing systems* 16. MIT Press
9. Drineas P, Kannan R, Mahoney M (2006) , Fast Monte Carlo algorithms III: computing a compressed approximate matrix decomposition. *SIAM J Comput* 36(1):184–206
10. Faloutsos C, Lin K-I (1995) FastMap: a fast algorithm for indexing, data-mining and visualization of traditional and multimedia datasets. In: *Proceedings of ACM SIGMOD conference*
11. Golub G, van Loan J (1996) *Matrix computations*. 3. Johns Hopkins University Press, Baltimore
12. Halevy A, Norvig P, Pereira F (2009) The unreasonable effectiveness of data. *IEEE Intell Syst* 24(2):8–12
13. Hoyer P (2004) Non-negative matrix factorization with sparseness constraints. *J Mach Learn* 5(Dec):1457–1469
14. Hueter I (1999) Limit theorems for the convex hull of random points in higher dimensions. *Trans Am Math Soc* 351(11):4337–4363
15. Jolliffe I (1986) *Principal component analysis*. Springer, New York
16. Kim J, Park H (2008) Toward faster nonnegative matrix factorization: a new algorithm and comparisons. In: *Proceedings of IEEE international conference on data mining*
17. Klingenberg B, Curry J, Dougherty A (2008) Non-negative matrix factorization: ill-posedness and a geometric algorithm. *Pattern Recogn* 42(5):918–928
18. Langville A, Meyer C, Albright R (2006) Initializations for the nonnegative matrix factorization. In: *Proceedings of ACM international conference on knowledge discovery and data mining*
19. Lee DD, Seung HS (1999) Learning the parts of objects by non-negative matrix factorization. *Nature* 401(6755):788–799
20. Li T (2008) Clustering based on matrix approximation: a unifying view. *Knowl Inf Syst* 17(1):1–15
21. Oliva A, Torralba A (2001) Modeling the shape of the scene: a holistic representation of the spatial envelope. *Int J Comput Vis* 42(3):145–175
22. Ostrouchov G, Samatova N (2005) On FastMap and the convex hull of multivariate data: toward fast and robust dimension reduction. *IEEE Trans Pattern Anal Mach Intell* 27(8):1340–1434
23. Paatero P, Tapper U (1994) Positive matrix factorization: a non-negative factor model with optimal utilization of error estimates of data values. *Environmetrics* 5(2):111–126
24. Rennie J, Srebro N (2005) Fast maximum margin matrix factorization for collaborative prediction. In: *Proceedings of international conference on machine learning*
25. Srebro N, Rennie JM, Jaakola T (2005) Maximum-margin matrix factorization. In: *Advances in neural information processing systems* 17. MIT Press
26. Sun J, Xie Y, Zhang H, Faloutsos C (2007) Less is more: compact matrix decomposition for large sparse graphs. In: *Proceedings of SIAM international conference on data mining*

27. Suvrit S (2008) Block-iterative algorithms for non-negative matrix approximation. In: Proceedings of IEEE international conference on data mining
28. Thureau C, Kersting K, Bauckhage C (2009) Convex non-negative matrix factorization in the Wild. In: Proceedings of IEEE international conference on data mining
29. Torralba A, Fergus R, Freeman WT (2008) 80 Million tiny images: a large data set for nonparametric object and scene recognition. *IEEE Trans Pattern Anal Mach Intell* 30(11):1958–1970
30. Vasiloglou N, Gray A, Anderson D (2009) Non-negative matrix factorization, convexity and isometry. In: Proceedings of SIAM international conference on data mining
31. Ziegler G (1995) Lectures on polytopes. Springer, New York

Author Biographies



Christian Thureau is a researcher at the Visual and Social Media Group at the Fraunhofer Institute for Intelligent Analysis and Information Systems (IAIS) and at the Bonn-Aachen International Center for Information Technology (B-IT). In 2006 he received a Ph.D. in computer science from Bielefeld University, Germany. From 2007 till 2008 he was a Marie Curie post-doctoral researcher at the Center for Machine Perception Prague, Czech Republic. In 2008 he stayed at the Technical University Dortmund, before joining Fraunhofer IAIS the same year. His main research interests are pattern recognition and data mining, and their application to new media data. He reviews for leading journals and conferences, and has organized several workshops and sessions in this area.



Kristian Kersting is the head of the STREAM group at Fraunhofer IAIS, Bonn, Germany, and a research fellow of the University of Bonn. He received his Ph.D. from the University of Freiburg, Germany, in 2006. After a PostDoc at the MIT, USA, he joined Fraunhofer IAIS in 2008 to build up the STREAM research group. His main research interests are statistical relational AI, machine learning, data mining, and robotics. He has published more than 60 papers, has received the ECML Best Student Paper Award in 2006 and the ECCAI Dissertation Award 2006 for the best European dissertation in the field of AI, and is an ERCIM Cor Baayen Award 2009 finalist. He served as area chair and PC member for numerous conferences in the area of machine learning and artificial intelligence. He is on the editorial board of the *Machine Learning Journal* and the *Journal of Artificial Intelligence Research*.



Mirwaes Wahabzada received his Diploma (M.Sc. equiv.) in computer science from the University of Bonn, Germany in 2009. He is currently a graduate student in the Knowledge Discovery Department at Fraunhofer IAIS in Bonn, Germany. His research interests include machine learning, graphical models, statistical relational learning, text/web mining.



Christian Bauckhage is professor of media informatics at the University of Bonn and heads the Visual and Social Media Group at the Fraunhofer Institute for Intelligent Analysis and Information Systems (IAIS). He obtained a Ph.D. in computer Science from Bielefeld University, Germany, and was a postdoctoral researcher at the Centre for Vision Research in Toronto, Canada. Later, he worked as a senior research scientist at Deutsche Telekom Laboratories in Berlin, where he conducted and coordinated industrial ICT research. He is a regular reviewer for leading journals and conferences and has (co)authored more than 100 technical papers on pattern recognition, web data mining, social media, and future Internet related topics.