# On the Relation Between Autoencoders and Non-negative Matrix Factorization, and Their Application for Mutational Signature Extraction

Ida Egendal[1,2], Rasmus Froberg Brøndum[1,2], Marta Pelizzola[3], Asger Hobolth[3], and Martin Bøgsted[1,2]

[1]Center for Clinical Data Science, Aalborg University and Aalborg University Hospital, Aalborg, Denmark
[2]Clinical Cancer Research Center, Aalborg University Hospital, Aalborg, Denmark
[3]Department of Mathematics, Aarhus University, Aarhus, Denmark

## Abstract

The aim of this study is to provide a foundation to understand the relationship between non-negative matrix factorization (NMF) and non-negative autoencoders enabling proper interpretation and understanding of autoencoder-based alternatives to NMF. Since its introduction, NMF has been a popular tool for extracting interpretable, low-dimensional representations of high-dimensional data. However, recently, several studies have proposed to replace NMF with autoencoders. This increasing popularity of autoencoders warrants an investigation on whether this replacement is in general valid and reasonable. Moreover, the exact relationship between non-negative autoencoders and NMF has not been thoroughly explored. Thus, a main aim of this study is to investigate in detail the relationship between non-negative autoencoders and NMF. We find that the connection between the two models can be established through convex NMF, which is a restricted case of NMF. In particular, convex NMF is a special case of an autoencoder. The performance of NMF and autoencoders is compared within the context of extraction of mutational signatures from cancer genomics data. We find that the reconstructions based on NMF are more accurate compared to autoencoders, while the signatures extracted using both methods show comparable consistencies and values when externally validated. These findings suggest that the non-negative autoencoders investigated in this article do not provide an improvement of NMF in the field of mutational signature extraction.

## 1 Introduction

Non-negative matrix factorization (NMF) is a popular tool for unsupervised learning [Lee and Seung, 1999]. NMF factorizes a non-negative data matrix into a product of two non-negative matrices of lower dimension: a basis matrix consisting of basis vectors and a weight matrix consisting of the basis vector's weights for each observation in the data matrix. NMF has gained a strong footing in different scientific fields due to its high interpretability [Alexandrov et al., 2013; Fang et al., 2018; Özer et al., 2022]. Specifically, NMF has proven to be a useful tool to derive mutational signatures from cancer genomics data.

In mutational signature analysis, it is typically assumed that all mutations in a cancer genome are caused by mutagenic processes that leave a characteristic pattern of mutations in the genome. These patterns are denoted mutational signatures. Several signatures have been identified and linked to different mutagenic processes such as ultraviolet light exposure and tobacco smoking [Nik-Zainal et al., 2015]. Alexandrov et al. (2013) proposed using NMF on mutational count data from cancer genomes to decipher the mutational signatures of the processes the patients have been exposed to throughout the development of the disease. NMF has since then been the dominating model for mutational signature extraction [Alexandrov et al., 2020; Blokzijl et al., 2018; Islam et al., 2022]. When extracting mutational signatures with NMF, the data matrix consisting of a number of patients' mutational profiles is decomposed into a matrix representing the signatures of mutagenic processes (basis vectors) and an exposure matrix

dictating the number of mutations that can be attributed to each specific process in the mutational profiles of each patient (weight matrix).

Recently, several studies have proposed substituting NMF with non-negative autoencoders which are increasingly popular for dimensionality reduction [Hosseini-Asl, Zurada, and Nasraoui, 2016; Khatib et al., 2018; Lemme, Reinhart, and Steil, 2012; Smaragdis and Venkataramani, 2017; Özer et al., 2022]. This is also the case in mutational signature extraction [Pancotti et al., 2023; Pei et al., 2020]. Pei et al. (2020) suggested using a sparse autoencoder to identify mutational signatures from cancer genomics data and generated estimates that were not only in concordance with existing literature but also correlated with observed exogenous exposures in a meaningful way. Pancotti et al. (2023) suggested a hybrid architecture with a deep encoding and shallow decoding to relax the assumption of linearity NMF imposes on mutational signature extraction. However, they did not compare the results and performance to NMF. This trend of using non-negative autoencoders as an alternative to NMF with promising results prompts the question: What is the mathematical relation between autoencoders and NMF? And how do they compare in applications?

The aim of this study is to compare the performance of shallow, non-negative autoencoders to NMF in the field of mutational signatures. In particular, we compare the in- and out-of-sample reconstruction error, as well as the consistency of the extracted signatures from the tumor-normal whole genome sequences of 713 ovary, 311 prostate, and 523 uterus tumors from the Genomics England 100,000 Genomes Project (GEL) [Turnbull, 2018; Turro et al., 2020]. Moreover, we show theoretically that shallow, non-negative autoencoders and NMF is a special case of NMF where the basis vectors are restricted to be convex combinations of columns in the data matrix [Ding, Li, and Jordan, 2010]. Based on this, we deduce how it impacts interpretation of the estimates generated by non-negative autoencoders in general and especially within the field of mutational signatures.

In Section 2, we characterize the mathematical relationship between autoencoders and NMF and introduce the framework for comparing the two models. In Section 3, we compare the performance of NMF and autoencoders and demonstrate the mathematical equivalence between convex NMF and non-negative, shallow autoencoders. Lastly, we discuss and conclude on the results in Sections 4 and 5.

## 2 Methods

Consider a non-negative data matrix $\boldsymbol{V} \in \mathbb{R}_+^{M \times N}$. In this study the aim is to decompose $\boldsymbol{V}$ into a basis matrix $\boldsymbol{H} \in \mathbb{R}_+^{M \times K}$ and a weight matrix $\boldsymbol{W} \in \mathbb{R}_+^{K \times N}$, where $M$ denotes the number of features, $N$ denotes the number of observations, and $K$ denotes the number of basis vectors in the latent representation of $\boldsymbol{V}$. A schematic overview of all considered decompositions in this section can be seen in Figure 1.

### 2.1 Non-negative matrix factorization

NMF decomposes a matrix with non-negative entries into a matrix product of two factor matrices with non-negative entries, one containing a set of basis vectors and one containing a set of weights. The shared dimension, $K$, of the factor matrices, is typically chosen to be much smaller than the dimensions of the input matrix, making NMF a dimensionality reduction technique.

Standard $K$-dimensional NMF was introduced by Lee and Seung (1999) and aims to make a reconstruction, $\hat{\boldsymbol{V}}$, of the original data matrix by a product of two non-negative matrices:

$$\hat{\boldsymbol{V}} = \boldsymbol{H}\boldsymbol{W}, \tag{1}$$

where each column in $\boldsymbol{H}$ represents a basis vector and each column in $\boldsymbol{W}$ represents each sample's weights when being reconstructed as a linear mixture of the basis vectors, i.e.,

$$\hat{\boldsymbol{v}}_n = \sum_{k=1}^{K} \boldsymbol{h}_k w_{k,n}, \quad n = 1, \ldots, N, \tag{2}$$

where $\hat{\boldsymbol{v}}_n$ is the $n$'th column of the reconstructed data matrix, $\hat{\boldsymbol{V}}$, $w_{n,k}$ represents the $(k,n)$'th entry of the weight matrix $\boldsymbol{W}$, and $\boldsymbol{h}_k$ represents the $k$'th column of the basis matrix $\boldsymbol{H}$ [Lee and Seung, 1999].
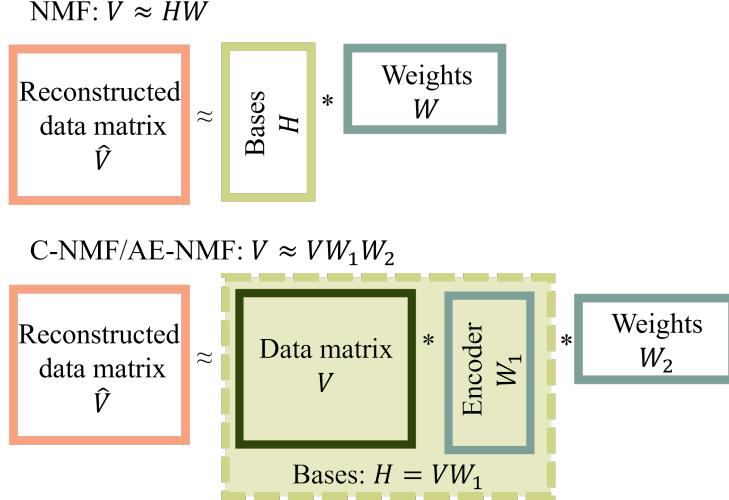
Figure 1: Schematic representation of the composition of basis vectors and weights in NMF (top) and C-NMF and AE-NMF (bottom).

## 2.2 Convex NMF

Convex NMF, as introduced by Ding, Li, and Jordan (2010), is a special case of NMF, where the basis vectors are constrained to be spanned by the columns of the data matrix $\boldsymbol{V}$, thus the data matrix is approximated by:

$$\hat{\boldsymbol{V}} = \boldsymbol{V}\boldsymbol{W}_1\boldsymbol{W}_2, \tag{3}$$

with $\boldsymbol{W}_1, \boldsymbol{W}_2^T \in \mathbb{R}_+^{N \times K}$. Defining $\boldsymbol{H} := \boldsymbol{V}\boldsymbol{W}_1$ and $\boldsymbol{W} := \boldsymbol{W}_2$ gives the model formulation of NMF as defined in Equation (1). Ding, Li, and Jordan (2010) focus mainly on the general case where the data matrix, $\boldsymbol{V}$, can assume values within all real numbers, but for this paper we consider a version of convex NMF where $\boldsymbol{V}$ is constrained to be non-negative.

## 2.3 Autoencoders

Autoencoders consist of an encoder applied to the input to create a latent representation and a decoder that maps the latent representation to a reconstruction of the input [Kramer, 1991]. Choosing the dimension of the latent representation to be lower than the dimension of the input makes the autoencoder a dimensionality reduction technique. A single hidden layer and fully connected autoencoder's reconstruction, $\hat{\boldsymbol{V}}$, of a data matrix, $\boldsymbol{V}$, is mathematically defined as:

$$\hat{\boldsymbol{V}} = \phi_{\text{dec}}(\phi_{\text{enc}}(\boldsymbol{V}\boldsymbol{W}_{\text{enc}} + \boldsymbol{b}_{\text{enc}})\boldsymbol{W}_{\text{dec}} + \boldsymbol{b}_{\text{dec}}), \tag{4}$$

where $\boldsymbol{W}_{\text{enc}}, \boldsymbol{W}_{\text{dec}}^T \in \mathbb{R}^{N \times K}$ are the encoding and decoding matrices, $\boldsymbol{b}_{\text{enc}} \in \mathbb{R}^K$ and $\boldsymbol{b}_{\text{dec}} \in \mathbb{R}^N$ are the bias terms of the encoding and the decoding layers and $\phi_{\text{enc}} : \mathbb{R}^{M \times K} \mapsto \mathbb{R}^{M \times K}$ and $\phi_{\text{dec}} : \mathbb{R}^{M \times N} \mapsto \mathbb{R}^{M \times N}$ are the activation functions which provide entry-wise modifications of the affected nodes.

## 2.4 Mathematical Equivalence and Interpretation

Setting $b_{\text{enc}} = \boldsymbol{0}_K$, $b_{\text{dec}} = \boldsymbol{0}_N$, and $\phi_{\text{enc}}, \phi_{\text{dec}} : x \mapsto x$ while constraining the weights to be non-negative in Equation (4) yields exactly the convex NMF formulation from Equation (3), where $\boldsymbol{W}_1$ corresponds to the encoding matrix and $\boldsymbol{W}_2$ corresponds to the decoding matrix. Thus, non-negative autoencoders can be constructed to be mathematically equivalent to convex NMF. The choice between convex NMF and the autoencoder defined above therefore reduces to a choice of how to optimize the problem, either by the multiplicative updating steps derived by Ding, Li, and Jordan (2010) or by the gradient descent-based additive updates of the autoencoder. The core architecture of this autoencoder is analogous to that of Pei et al. (2020), but differ by fixing $b_{\text{enc}}$ and $b_{\text{dec}}$ to zero instead of estimating them through training and

choosing the identity function as activation functions instead of $\phi_{\text{enc}} : A \mapsto ReLU(x) = \max(0, A)$ and $\phi_{\text{dec}} : A \mapsto Softmax(A) = \{\exp(a_{m,n})/\sum_{m=1}^{M} \exp(a_{m,n})\}_{m=1,\ldots,M,n=1,\ldots,N}$, for a matrix $\boldsymbol{A} \in \mathbb{R}^{M \times N}$.

We will use 'C-NMF' to refer to convex NMF optimized with the multiplicative updating steps derived by Ding, Li, and Jordan (2010) and use 'AE-NMF' for the class of autoencoders constructed equivalently to convex NMF.

The encoded data matrix in AE-NMF and C-NMF, $\boldsymbol{VW}_1$, which is a convex combination of the columns of the data matrix, is interpreted as the basis matrix, $\boldsymbol{H}$, in conventional NMF and hereby the interpretation of $\boldsymbol{W}_2$ corresponds to that of $\boldsymbol{W}$ in conventional NMF. Specifically, within mutational signature analysis this means that the signature matrix is a convex combination of the patients' mutational profiles which aligns with how mutational signatures are understood. If the data matrix is transposed to follow conventional orientation in neural networks, i.e., with patients as rows and mutation types as columns the weights or exposures are modelled as convex combinations of the mutations which lacks an equally straightforward interpretation and direct equivalence with convex NMF. The equivalence between C-NMF and AE-NMF is the necessary and previously missing link that enables one to interpret the parameters in AE-NMF similarly to NMF, therefore this orientation is crucial for proper comparison.

Though the interpretation of AE-NMF, C-NMF, and NMF is similar, there are still considerable differences between C-NMF, AE-NMF, and standard NMF. In particular, NMF estimates $N \cdot K + K \cdot M$ parameters whereas AE-NMF and C-NMF estimates $2 \cdot (K \cdot N)$ parameters. Thus, AE-NMF and C-NMF will estimate a larger number of parameters in the factor matrices than NMF when the number of observations $N$ surpasses the number of features $M$, which is often the case within mutational signatures.

## 2.5   Modeling and performance framework

All estimation in this paper, for NMF, C-NMF, and AE-NMF, is done by minimizing the average Frobenius distance between the original input and the reconstructed input:

$$L_F(\boldsymbol{V}, \hat{\boldsymbol{V}}) = \frac{||\boldsymbol{V} - \hat{\boldsymbol{V}}||_F}{M \cdot N} = \frac{1}{M \cdot N} \sqrt{\sum_{i=1}^{M} \sum_{j=1}^{N} |v_{i,j} - \hat{v}_{i,j}|^2}. \tag{5}$$

NMF is estimated using the multiplicative and iterative updating steps defined by Lee and Seung (2000):

$$\boldsymbol{H} \leftarrow \boldsymbol{H} \frac{\boldsymbol{VW}^T}{\boldsymbol{HWW}^T}, \qquad \boldsymbol{W} \leftarrow \boldsymbol{W} \frac{\boldsymbol{H}^T\boldsymbol{V}}{\boldsymbol{H}^T\boldsymbol{HW}}. \tag{6}$$

C-NMF is estimated using the multiplicative updating scheme defined by Ding, Li, and Jordan (2010). This updating scheme reduces to

$$\boldsymbol{W}_1 \leftarrow \boldsymbol{W}_1 \sqrt{\frac{(\boldsymbol{V}^T\boldsymbol{V})\boldsymbol{W}_2^T}{(\boldsymbol{V}^T\boldsymbol{V})\boldsymbol{W}_1\boldsymbol{W}_2\boldsymbol{W}_2^T}}, \qquad \boldsymbol{W}_2^T \leftarrow \boldsymbol{W}_2^T \sqrt{\frac{(\boldsymbol{V}^T\boldsymbol{V})\boldsymbol{W}_1}{\boldsymbol{W}_2^T\boldsymbol{W}_1^T(\boldsymbol{V}^T\boldsymbol{V})\boldsymbol{W}_1}}, \tag{7}$$

in the case where $\boldsymbol{V}$ is constrained to be in the non-negative domain.

AE-NMF is estimated by an Adam optimizer, where the $t$'th iteration of a parameter, $w^{(t)}$, is updated additively and iteratively using the following method:

$$w^{(t+1)} \leftarrow w^{(t)} - \eta \frac{\frac{m^{(t)}}{1-\beta_1^t}}{\sqrt{\frac{v^{(t)}}{1-\beta_2^t}} + \varepsilon}, \quad \text{where} \quad m^{(t)} \leftarrow \beta_1 m^{(t-1)} + (1-\beta_1)\nabla_w L_t$$
$$\text{and} \quad v^{(t)} \leftarrow \beta_2 v^{(t-1)} + (1-\beta_2)(\nabla_w L_t)^2. \tag{8}$$

Here $L_t$ is the loss function at the $t$'th iteration and $\varepsilon = 10^{-8}$ is a small scalar to prevent division by zero. The factors $\beta_1$ and $\beta_2$ are exponential decay rates initialized to the default values $\beta_1 = 0.9$ and $\beta_2 = 0.999$. The moment vectors of parameter $w$ at time $t$, $m^{(t)}$ and $v^{(t)}$, are initialized as $m^{(0)} \leftarrow 0$ and $v(0) \leftarrow 0$ [Kingma and Ba, 2017]. Training of AE-NMF is done using the PyTorch module, version 1.13.1 [Paszke et al., 2019], in Python version 3.11.3 [Van Rossum and Drake, 2009].

### 2.5.1 Average Cosine Similarity

The similarity between two sets of basis vectors is evaluated using the average cosine similarity (ACS). The cosine similarity for two equal length vectors, $\tilde{\boldsymbol{h}}, \hat{\boldsymbol{h}} \in \mathbb{R}^M$, is defined as

$$S_C(\tilde{\boldsymbol{h}}, \hat{\boldsymbol{h}}) = \frac{\tilde{\boldsymbol{h}} \cdot \hat{\boldsymbol{h}}}{||\tilde{\boldsymbol{h}}|| ||\hat{\boldsymbol{h}}||}. \tag{9}$$

The ACS for two matrices of matched basis vectors $\tilde{\boldsymbol{H}}, \hat{\boldsymbol{H}} \in \mathbb{R}_+^{M \times K}$ with $K$ basis vectors is defined as

$$ACS(\tilde{\boldsymbol{H}}, \hat{\boldsymbol{H}}) = \frac{1}{K} \sum_{k=1}^{K} S_c(\tilde{\boldsymbol{h}}_k, \hat{\boldsymbol{h}}_k), \tag{10}$$

where $\tilde{\boldsymbol{h}}_k$ denotes the $k$'th column in $\tilde{\boldsymbol{H}}$, and $\hat{\boldsymbol{h}}_k$ is the corresponding basis vector in $\hat{\boldsymbol{H}}$.

### 2.5.2 Signature Matching

Given two matrices of basis vectors $\boldsymbol{H} \in \mathbb{R}^{M \times K}$ and $\tilde{\boldsymbol{H}} \in \mathbb{R}^{M \times \tilde{K}}$ where $K \leq \tilde{K}$, the task is to find $K$ pairs $(i, j)$, where $i = 1, \ldots, K$ and $j \in \{1, \ldots, \tilde{K}\}$, such that the sum of cosine similarities (and thus the ACS) over all pairs $\sum_{(i,j)} S_c(\boldsymbol{h}_i, \tilde{\boldsymbol{h}}_j)$ is maximized, all vectors in $\boldsymbol{H}$ are matched exactly once, and all vectors in $\tilde{\boldsymbol{H}}$ are matched at most once. This combinatorial problem is a linear assignment problem, where the cost matrix is all cosine distances between two basis vectors $\left\{1 - S_c(\boldsymbol{h}_i, \tilde{\boldsymbol{h}}_j)\right\}_{i=1,\ldots,K, j=1,\ldots,\tilde{K}}$. This problem is solved using the Hungarian algorithm in both the balanced ($K = \tilde{K}$) and unbalanced ($K < \tilde{K}$) case [Kuhn, 1955].

## 3 Results

In this section, we compare NMF, C-NMF, and AE-NMF in a simple simulated example (Section 3.1) and compare NMF and AE-NMF on the ovary, prostate, and uterus genomes from the Genomics England 100,000 Genomes Project [Turnbull, 2018; Turro et al., 2020] by the ability to reconstruct the input data accurately and to generate stable and sensible mutational signatures (Section 3.2).

All training was performed with a relative tolerance of $10^{-10}$ as convergence criteria. Estimation with AE-NMF was performed with an Adam optimizer with a learning rate of $10^{-4}$. Non-negativity in AE-NMF was enforced taking the absolute value of the encoding and decoding weight matrices in the forward pass. The choice of method to enforce non-negativity in AE-NMF is elaborated in Supplementary Material Section 1.4.

### 3.1 Simulated example

Consider an example with two basis vectors (signatures) consisting of 6 features (mutation types):

$$\begin{aligned} \boldsymbol{h}_1 &= (2, 2, 1, 1, 0, 0)^T / 6 \\ \boldsymbol{h}_2 &= (0, 0, 0, 1, 1, 1)^T / 3. \end{aligned} \tag{11}$$

From these basis vectors we simulated 30 samples (patients), $\{\boldsymbol{v}_1, \ldots, \boldsymbol{v}_{30}\}$, where each sample was simulated with one of three distinct compositions of $\boldsymbol{h}_1$ and $\boldsymbol{h}_2$ and Poisson distributed noise:

$$\boldsymbol{v}_j = (v_{j,1}, v_{j,2}, v_{j,3}, v_{j,4}, v_{j,5}, v_{j,6})^T, \tag{12}$$

for $j = 1, \ldots, 30$ and

$$v_{i,j} \overset{\text{i.i.d}}{\sim} \begin{cases} \text{Po}(180 \cdot \boldsymbol{h}_{i,1} + 20 \cdot \boldsymbol{h}_{i,2}) & \text{for } j = 1, \ldots, 10 \\ \text{Po}(100 \cdot \boldsymbol{h}_{i,1} + 100 \cdot \boldsymbol{h}_{i,2}) & \text{for } j = 11, \ldots, 20 \\ \text{Po}(20 \cdot \boldsymbol{h}_{i,1} + 180 \cdot \boldsymbol{h}_{i,2}) & \text{for } j = 21, \ldots, 30, \end{cases} \tag{13}$$

for $i = 1, \ldots, 6$. For the full data matrix, $\boldsymbol{V} = [v_{i,j}]$, the two basis vectors and the corresponding weights (exposures) were estimated using NMF, C-NMF, and AE-NMF, with the Frobenius norm as the loss
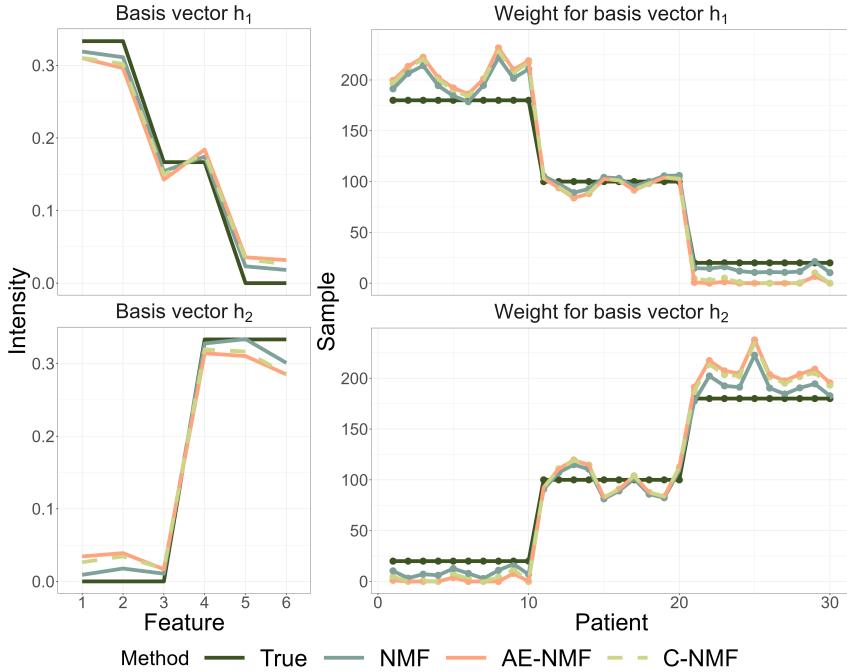
Figure 2: Estimated and true weights and basis vectors of the simulated data. Note the light green lines representing C-NMF coincide with the the dashed, coral lines of AE-NMF since the estimates found by both methods are almost identical.

function. To encourage convergence towards the same minimum, C-NMF and AE-NMF were initialized with the same matrices, $\boldsymbol{W}_1^{(0)}$ and $\boldsymbol{W}_2^{(0)}$, respectively to a $30 \times 2$ and a $2 \times 30$ matrix with values sampled uniformly in $[0, 1)$. This was not possible for NMF since its factor matrices have different dimensions. For NMF the matrix of basis vectors was initialized uniformly in $[0, 1)$, and the weight matrix was initialized as $\boldsymbol{W}_2^{(0)}$ from C-NMF and AE-NMF.

Figure 2 shows the estimated basis vectors and weights for each method. In this example, identical factor matrices are recovered by the C-NMF and AE-NMF updates, illustrated by the coinciding dashed green and coral lines, while NMF yields different results, illustrated by the blue lines. Furthermore, the NMF solution reconstructs the input with higher accuracy, with a reconstruction error of 12.95 compared to the higher values of 13.64 for AE-NMF and 13.00 for C-NMF.

## 3.2 Cancer Data Analysis

To emulate the terminology used in the cancer data analysis, the matrix of basis vectors will be denoted as the signatures, and the weight matrix will be denoted the exposures.

The number of signatures in each diagnosis was determined with the method described in the Supplementary Material Section 1.3 with ten training/test splits for each diagnosis. The test errors as a function of $K \in \{2, \ldots, 12\}$ are depicted in Figure S1. Algorithm S2 in the Supplementary Material yielded three signatures for AE-NMF and four signatures for C-NMF and NMF in the ovary cohort; four, four and six signatures for AE-NMF, C-NMF and NMF, respectively in the prostate cohort and four, six, and 11 signatures for AE-NMF, C-NMF, and NMF respectively in the uterus cohort. Thus, the number of signatures used in the cancer data analyses is chosen, using the weighted average in Equation (S4) in the Supplementary Material, as four for the ovary cohort, five for the prostate cohort and eight for the uterus cohort.

As shown in Section 2.4 C-NMF and AE-NMF are mathematically equivalent. Furthermore, since the resulting factor matrices of C-NMF and AE-NMF were identical in the simulated example and performed similarly in the cancer data analysis with respect to reconstruction error (Figure S1 and Figure S3) and consistency (Figure S4), we consider C-NMF and AE-NMF as practically equivalent. Thus, the following

| Cohort | Split | $n$ | Average Error | | |
| --- | --- | --- | --- | --- | --- |
| | | | NMF | AE-NMF | Ratio |
| Ovary | Train | 418 | $\mathbf{5.92 \cdot 10^3}$ | $1.29 \cdot 10^4$ | 2.17 |
| | Test | 105 | $\mathbf{1.51 \cdot 10^6}$ | $1.72 \cdot 10^6$ | 1.63 |
| Prostate | Train | 248 | $\mathbf{1.06 \cdot 10^2}$ | $2.39 \cdot 10^2$ | 2.24 |
| | Test | 63 | $\mathbf{2.81 \cdot 10^3}$ | $3.16 \cdot 10^3$ | 1.23 |
| Uterus | Train | 570 | $\mathbf{5.33 \cdot 10^5}$ | $1.39 \cdot 10^6$ | 2.61 |
| | Test | 143 | $\mathbf{1.11 \cdot 10^6}$ | $2.57 \cdot 10^6$ | 2.45 |

Table 1: Average training and test error between the input and reconstructed data for each method across the 30 splits of the ovary, prostate, and uterus cohort and the average ratio between the errors. The lowest reconstruction error in each cohort is highlighted in bold.

analyses will be performed comparing only NMF and AE-NMF.

### 3.2.1 Extraction performance

For each cohort, we divided the patients' mutational profile data into 30 80/20 train/test set splits. *De novo* extractions by NMF and AE-NMF were performed on the training set yielding a set of training errors and refits were performed on the test sets yielding a set of test errors. Plots of the first and second principal component of all extracted signatures from the 30 training sets are depicted in Figure 3.

The procedure of calculating the test errors is detailed in the Supplementary Material Section 1.2. Boxplots of the training and test errors when reconstructing the input matrix for each method and diagnosis are shown in Figure 4. Table 1 reports the average training and test error across the 30 splits for each method and diagnosis. Considering the reconstruction errors in Figure 4 and Table 1, it is apparent that NMF consistently performs better than AE-NMF in terms of reconstructing the input data. This is the case both on average and in the majority of splits as shown by Figure 4. The ratios in Table 1 reveal that the difference is more expressed in the training splits than the test splits.

For each method and diagnosis the analyses yielded 30 signature sets, one from each training set. The consistency of the estimated signatures extracted within each method is investigated by calculating the ACS between each pairwise combination of signatures extracted using a given method across the 30 splits of the data matrix. This yields a total of $\binom{30}{2} = 435$ comparisons for each cohort and method, and resulted in an average ACS consistency of 0.91 for NMF and 0.85 for AE-NMF in the ovary cohort, 0.86 for NMF and 0.88 for AE-NMF in the prostate cohort, and 0.94 for NMF and 0.91 for AE-NMF in the uterus cohort. To asses whether the difference in mean consistency between NMF and AE-NMF were significant, two-sided t-tests for equal means were performed for each cohort. These resulted in a $p$-value of $3.7 \cdot 10^{-18}$ for the ovary cohort, $7.0 \cdot 10^{-10}$ for the prostate cohort, and $1.1 \cdot 10^{-11}$ for the uterus cohort, thus the mean consistency for NMF and AE-NMF can not be assumed to be equal for any cohort. The consistencies are depicted by boxplots in Figure 5, with green diamonds marking the average consistencies listed above. This reveals generally comparable consistencies of NMF and AE-NMF signatures.

### 3.2.2 COSMIC Validation

To compare the estimated signatures to those of the leading library of mutational signatures, COSMIC v. 3.4 [Tate et al., 2018], the signatures were clustered to form sets of consensus signatures using the partitioning around medoids (PAM) algorithm [Kaufman and Rousseeuw, 1990] with the number of clusters equal to the number of signatures used in the initial extractions. The clustering is depicted in Figure 3 where the first and second principal components of all *de novo* signatures are depicted and points are colored by their assigned PAM clustering. The consensus signatures were subsequently matched to the COSMIC signatures. The matched COSMIC signatures along with their cosine similarity can be seen in Table 2.

The ACS between the consensus signatures and their matched COSMIC signatures is 0.80 for NMF and 0.82 for AE-NMF in the ovary cohort, 0.90 for NMF and 0.86 for AE-NMF in the prostate cohort and 0.90 for NMF and 0.83 for AE-NMF in the uterus cohort. This reveals a slightly higher average similarity to COSMIC for NMF than AE-NMF. Additionally, the majority of NMF signatures also have better COSMIC matches than the corresponding AE-NMF signature. The matched COSMIC signatures
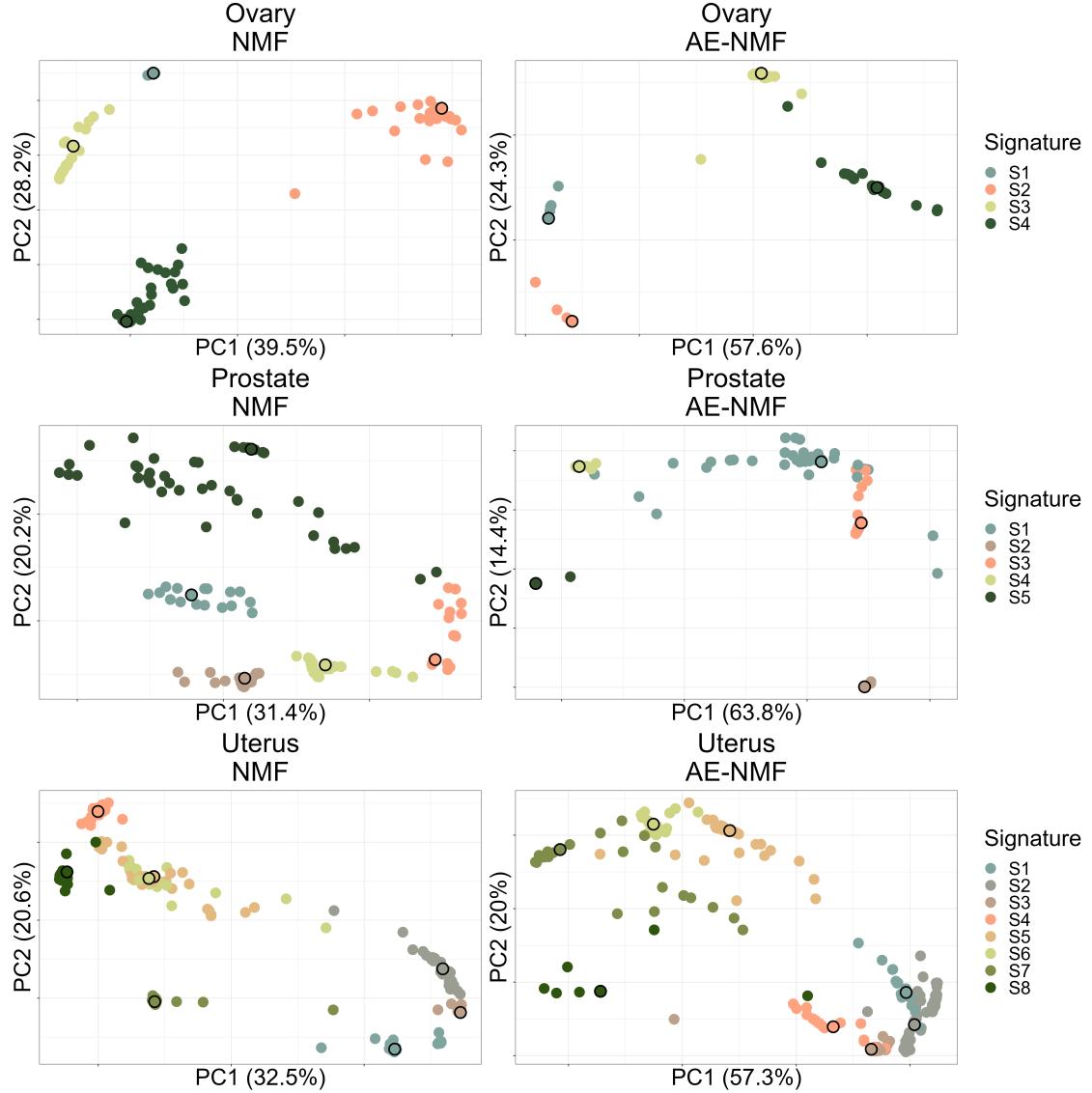
Figure 3: Second principal component plotted agains the first principal component of the *de novo* extracted signatures from the 30 train/test splits for AE-NMF and NMF (columns) and each diagnosis (rows). Points are colored by the PAM clustering assignment, and the cluster mediod is highlighted with a black outline. NMF: non-negative matrix factorization; AE: autoencoder; PAM: partition around mediods.
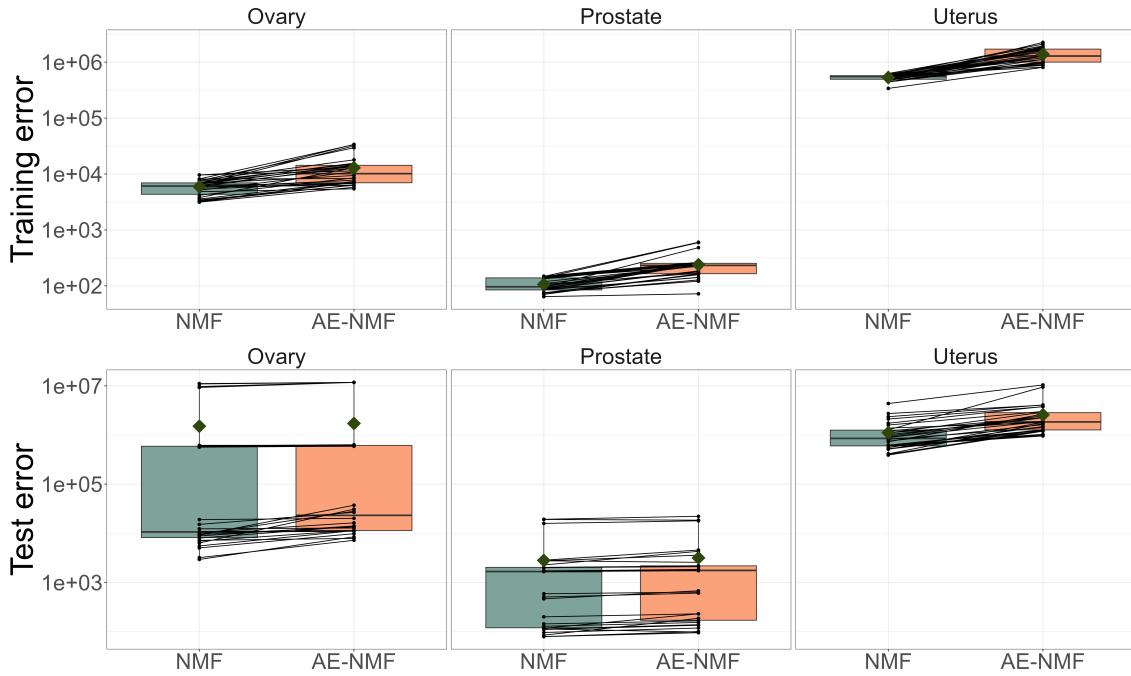
Figure 4: Boxplots of the training and test errors of 30 train/test splits of the ovary, prostate and uterus cohorts. NMF and AE-NMF errors resulting from the same splits are connected by a black line. The boxes are colored corresponding to the method used, and a green diamond depicts the average error. The y-axis is on log 10 scale.
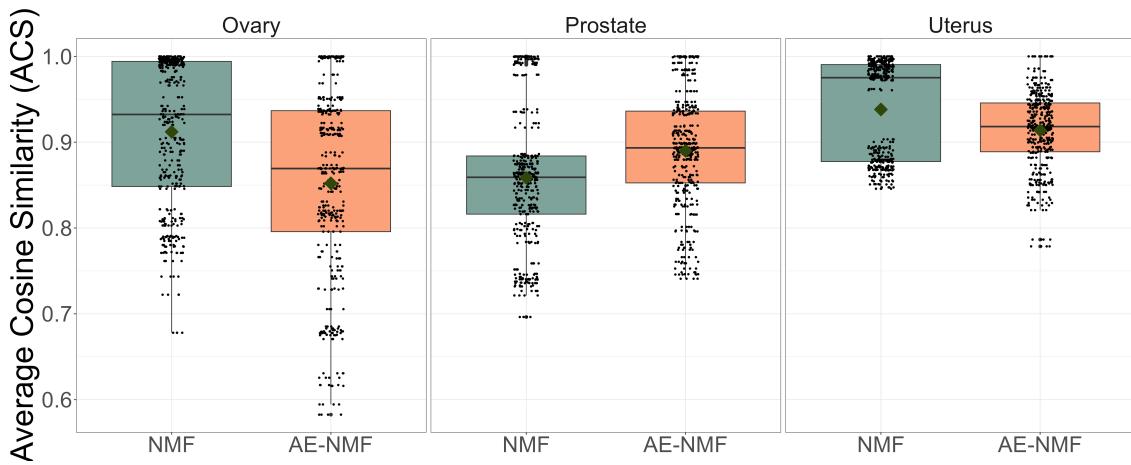


Figure 5: Boxplots of the ACS between each combination of signatures, extracted using NMF or AE-NMF and across the 30 splits of the data matrix for the ovary, prostate, and uterus cohort. The average is marked by a green diamond.

| Cohort | Signature | NMF | | AE-NMF | |
|---|---|---|---|---|---|
| **Ovary** | S1 | SBS10a | 0.93 | SBS10a | 0.93 |
| | S2 | **SBS3** | 0.95 | SBS40a | 0.80 |
| | S3 | SBS10c | 0.70 | SBS10c | 0.74 |
| | S4 | SBS44 | 0.84 | SBS44 | 0.82 |
| | ACS | | 0.80 | | 0.82 |
| **Prostate** | S1 | SBS17b | 0.84 | SBS17b | 0.77 |
| | S2 | **SBS33** | 0.95 | **SBS33** | 0.93 |
| | S3 | **SBS1** | 0.94 | **SBS6** | 0.90 |
| | S4 | SBS44 | 0.86 | SBS44 | 0.85 |
| | S5 | SBS40a | 0.90 | SBS40a | 0.87 |
| | ACS | | 0.90 | | 0.86 |
| **Uterus** | S1 | SBS10d | 0.96 | SBS10d | 0.95 |
| | S2 | SBS10c | 0.84 | **SBS10a** | 0.95 |
| | S3 | **SBS10a** | 0.97 | SBS10c | 0.75 |
| | S4 | **SBS6** | 0.88 | **SBS6** | 0.84 |
| | S5 | **SBS20** | 0.93 | **SBS20** | 0.90 |
| | S6 | **SBS10b** | 0.78 | **SBS2** | 0.71 |
| | S7 | **SBS28** | 0.96 | **SBS28** | 0.70 |
| | S8 | **SBS26** | 0.83 | **SBS26** | 0.82 |
| | ACS | | 0.89 | | 0.83 |

Table 2: The COSMIC signatures matched to the consensus signatures and the corresponding cosine similarity. The average cosine similarity (ACS) between a consensus signature set and the matched COSMIC signatures is reported in the last line for each diagnosis. Signatures that have previously been observed in patients with the same diagnosis are highlighted with bold text.

for all splits before clustering for NMF and AE-NMF can be seen in Supplementary Figure S5-S7. Of all matched consensus signatures, the following proportion has been previously observed in the corresponding diagnosis for each cohort: 1/4 of NMF and 0/4 of AE-NMF signatures in the ovary cohort, 2/5 of both NMF and AE-NMF signatures in the prostate cohort, and 6/8 of both NMF and AE-NMF signatures in the uterus cohort.

Overall both NMF and AE-NMF show a high degree of conformity with the COSMIC SBS signatures, and the extracted signatures are relevant to the diagnoses in which they have been identified. In signature cosistency, conformity with COSMIC and choosing relevant signatures, NMF and AE-NMF perform similarly, perhaps with a slight advantage to NMF.

# 4 Discussion

In this study, we compare NMF and AE-NMF by their ability to extract valid and consistent basis vectors and creating accurate reconstructions of the input data. We assert that such comparisons are theoretically meaningful since we demonstrate that AE-NMF and C-NMF are mathematically equivalent.

The study focuses on extracting mutational signatures in the ovary, prostate, and uterus cancer genomes of Genomics England's 100,000 Genomes cohort. NMF consistently outperformed AE-NMF in terms of reconstruction error; the differences being more expressed in the training splits than in the test splits. While AE-NMF constrains parameters to convex combinations of patients' profiles, NMF can freely assume non-negative values, giving it an advantage in training set reconstruction. When reconstructing the test set the signature matrix is fixed, and the task is thus identical for NMF and AE-NMF. One could expect the constrained nature of AE-NMF to regularize the signatures such that they reconstruct the test splits better compared to NMF, but as this was not the case in this study, it suggests that AE-NMF signatures may be generally less informative than the corresponding NMF signatures. The COSMIC validation revealed that both models recovered relevant signatures with high cosine similarity, with a slight advantage to NMF. This advantage is likely driven by the fact that the majority of COSMIC signatures are extracted using NMF-based methods [Alexandrov et al., 2020].

The mathematical equivalence between convex NMF and non-negative autoencoders enables the interpretation of parameters from AE-NMF to be similar to that of NMF. Thus, the mathematical

equivalence between convex NMF and non-negative autoencoders identified in this study is a necessary link in properly comparing AE-NMF and NMF. A link that has been missing in previous attempts to replace NMF with autoencoders while using the same interpretation. Squires, Bennett, and Niranjan (2019) also came to the conclusion that an autoencoder with the architecture of AE-NMF yields a hidden layer consisting of convex combinations of the data points but did not make the connection to convex NMF.

In practice we observed that C-NMF and AE-NMF will yield identical solutions in a sufficiently simple setup. When increasing the complexity of the problem both methods perform similarly in terms of reconstruction error and stability of the basis vectors, albeit finding different solutions within this minimum.

The architecture of AE-NMF is atypical for autoencoders by its shallow and linear nature and by transposing the input data matrix. By orienting the input matrix, $\boldsymbol{V}$, with features ($M$) as rows and observations ($N$) as columns the architecture will stray from how data is conventionally passed through neural networks, but this orientation is not uncommon in the literature of non-negative autoencoders [Khatib et al., 2018; Pei et al., 2020; Squires, Bennett, and Niranjan, 2019]. One direct consequence of this orientation is that it will not be possible to fit a separate sample by passing it through the trained network, a compelling attribute of neural networks, as the trained parameters will be observation-specific. Furthermore, the shallowness of AE-NMF favours the capture of linear relationships over more complex patterns. It is these strict architectural choices that yields the equivalence to convex NMF, ensuring a precise understanding of the interpretation of the parameters. Other efforts, such as the MUSE-XAE autoencoder for mutational signature extraction [Pancotti et al., 2023], have utilized the benefits of classic autoencoders with deeper extraction and conventional orientation of the input data matrix but at the cost of jeopardizing the link to NMF and thus the exact interpretation of the parameters, since this autoencoder is not equivalent to C-NMF.

All methods were optimized by minimizing the Frobenius loss function since this is the only loss function for conventional convex NMF by Ding, Li, and Jordan (2010) that presently has updating steps. This loss is only asymptotically efficient given Gaussian distributed input data, which is a questionable assumption since the data considered in this study consists of count data. Choosing a loss function adapted to Poisson may therefore be more adequate for this problem [Alexandrov et al., 2020; Degasperi et al., 2022]. An even more sophisticated error model would be the Negative Binomial distribution [Pelizzola, Laursen, and Hobolth, 2023]. The Negative Binomial distribution can model overdispersion in the mutational counts, which is often present on the patient-specific level. In future work, it would be interesting to see how training with a more appropriate loss function affects how NMF and AE-NMF compare in reconstruction accuracy and stability and whether the results from this study generalize to the Kullback-Leibler divergence.

In this study a relative convergence criteria of $10^{-10}$ was used for all analyses. Using such a low convergence criteria promoted that the three methods, based on two vastly different updating schemes, converged towards similar minima and, thus, established a common ground for comparison. A lower relative tolerance in the real data analyses was computationally infeasible. In particular, the high tolerance made especially C-NMF extremely time consuming in cases with many patients and/or signatures. The bootstrap method for determining the number of signatures were limited to ten splits for the same reason. If the aim is to just extract signatures using either method, we do not necessarily recommend training with such a low tolerance.

## 5 Conclusion

This study compares NMF with non-negative autoencoders, facilitated by the mathematical equivalence between non-negative autoencoders and convex NMF. This bridges a crucial gap in the comparison of NMF and non-negative autoencoders by offering insights into parameter interpretation that were previously lacking. The choice between convex NMF and its autoencoder equivalent is a question of choosing between a multiplicative or gradient descent-based optimizing algorithm to solve the same optimization problem. Thus, the non-negative autoencoder described in this study can be used as a faster alternative to convex NMF. Non-negative autoencoders exhibit higher reconstruction errors and similar consistencies compared to NMF, therefore the non-negative autoencoder investigated in this study is not a suitable alternative to NMF in mutational signature extraction. Thus, this study underscores the significance of methodological considerations when replacing NMF with non-negative autoencoders.

On the other hand autoencoders hold promise for modeling non-linearity, a capability absent in NMF, but such advancements are made at the cost of exact parameter interpretation.

# Acknowledgements

# Code- and Data avaliability

Code used for this study can be found at Github and the Genomics England WGS data used in this study was provided by Turro et al. (2020) on Zendo.

# Supplementary Materials

## Supplementary Methods

The supplementary methods contains a section on how test errors are calculated for all methods, the algorithms for choosing the optimal number of basis vectors, and the theoretical overview of methods for enforcing non-negativity in the autoencoder.

## Supplementary Results

The supplementary results contains the results used to choose the optimal number of basis vectors, the results used to determine the method to constrain non-negativity, and additional figures and tables.

# References

Alexandrov, Ludmil et al. (2013). "Deciphering Signatures of Mutational Processes Operative in Human Cancer". In: *Cell Reports* 3, pp. 246–259.

Alexandrov, Ludmil et al. (2020). "The repertoire of mutational signatures in human cancer". In: *Nature* 578, pp. 94–101.

Blokzijl, Francis et al. (2018). "MutationalPatterns: Comprehensive genome-wide analysis of mutational processes". In: *Genome Medicine* 10.

Degasperi, Andrea et al. (2022). "Substitution mutational signatures in whole-genome–sequenced cancers in the UK population". In: *Science* 376.

Ding, Chris H.Q., Tao Li, and Michael I. Jordan (2010). "Convex and Semi-Nonnegative Matrix Factorizations". In: *IEEE Transactions on Pattern Analysis and Machine Intelligence* 32, pp. 45–55.

Fang, Hao et al. (2018). "Sparsity-Constrained Deep Nonnegative Matrix Factorization for Hyperspectral Unmixing". In: *IEEE Geoscience and Remote Sensing Letters* 15, pp. 1105–1109.

Hosseini-Asl, Ehsan, Jacek Zurada, and Olfa Nasraoui (2016). "Deep Learning of Part-Based Representation of Data Using Sparse Autoencoders With Nonnegativity Constraints". In: *IEEE Transactions on Neural Networks and Learning Systems* 27, pp. 2486–2498.

Islam, S.M. Ashiqul et al. (2022). "Uncovering novel mutational signatures by de novo extraction with SigProfilerExtractor". In: *Cell Genomics* 2, p. 100179.

Kaufman, Leonard and Peter J. Rousseeuw (1990). "Partitioning Around Medoids (Program PAM)". In: *Finding Groups in Data*. John Wiley & Sons, Ltd. Chap. 2, pp. 68–125. eprint: https://onlinelibrary.wiley.com/doi/pdf/10.1002/9780470316801.ch2.

Khatib, Alaa El et al. (2018). "Nonnegative Matrix Factorization Using Autoencoders And Exponentiated Gradient Descent". In: *Int. Joint Conf. on Neural Networks (IJCNN)*, pp. 1–8.

Kingma, Diederik P. and Jimmy Ba (2017). "Adam: A Method for Stochastic Optimization". In: *arXiv preprint arXiv:1412.6980*.

Kramer, Mark A. (1991). "Nonlinear principal component analysis using autoassociative neural networks". In: *AIChE Journal* 37, pp. 233–243.

Kuhn, Harold (1955). "The Hungarian Method for The Assignment Problem". In: *Naval Research Logistics Quarterly* 2, pp. 83 –97.

Lee, Daniel and H. Sebastian Seung (2000). "Algorithms for Non-negative Matrix Factorization". In: *Advances in Neural Information Processing Systems*. MIT Press, pp. 556–562.

Lee, Daniel and Hyunjune Seung (1999). "Learning the Parts of Objects by Non-Negative Matrix Factorization". In: *Nature* 401, pp. 788–791.

Lemme, Andre, René Felix Reinhart, and Jochen Jakob Steil (2012). "Online learning and generalization of parts-based image representations by non-negative sparse autoencoders". In: *Neural Networks* 33, pp. 194–203.

Nik-Zainal, Serena et al. (2015). "The genome as a record of environmental exposure". In: *Mutagenesis* 30.

Pancotti, Corrado et al. (2023). "MUSE-XAE: MUtational Signature Extraction with eXplainable AutoEncoder enhances tumour type classification". In: *bioRxiv preprint 2023.10.23.562664*.

Paszke, Adam et al. (2019). "PyTorch: An Imperative Style, High-Performance Deep Learning Library". In: *Advances in Neural Information Processing Systems*.

Pei, Guangsheng et al. (2020). "Decoding whole-genome mutational signatures in 37 human pan-cancers by denoising sparse autoencoder neural network". In: *Oncogene* 39, 5031–5041.

Pelizzola, Marta, Ragnhild Laursen, and Asger Hobolth (2023). "Model selection and robust inference of mutational signatures using Negative Binomial non-negative matrix factorization". In: *BMC Bioinformatics* 24.

Smaragdis, Paris and Shrikant Venkataramani (2017). "A neural network alternative to non-negative audio models". In: *IEEE Int. Conf. on Acoustics, Speech and Signal Processing (ICASSP)*, pp. 86–90.

Squires, Steven, Adam Prügel Bennett, and Mahesan Niranjan (2019). "A Variational Autoencoder for Probabilistic Non-Negative Matrix Factorisation". In: *arXiv preprint arXiv:1906.05912*.

Tate, John G et al. (2018). "COSMIC: the Catalogue Of Somatic Mutations In Cancer". In: *Nucleic Acids Research* 47.D1, pp. D941–D947.

Turnbull, Clare (2018). "Introducing Whole Genome Sequencing into routine cancer care: The Genomics England 100,000 Genomes project". In: *Annals of Oncology* 29, pp. 784–787.

Turro, Ernest et al. (2020). "Whole-genome sequencing of patients with rare diseases in a national health system". In: *Nature* 583, 96–102.

Van Rossum, Guido and Fred L. Drake (2009). *Python 3 Reference Manual*. Scotts Valley, CA: CreateSpace.

Özer, Yigitcan et al. (2022). "Investigating Nonnegative Autoencoders for Efficient Audio Decomposition". In: *30th European Signal Processing Conf. (EUSIPCO)*, pp. 254–258.

# Supplementary Material for On the Relation Between Autoencoders and Non-negative Matrix Factorization, and Their Application for Mutational Signature Extraction

Ida Egendal[1,2], Rasmus Froberg Brøndum[1,2], Marta Pelizzola[3], Asger Hobolth[3], and Martin Bøgsted[1,2]

[1]Center for Clinical Data Science, Aalborg University and Aalborg University Hospital, Aalborg, Denmark
[2]Clinical Cancer Research Center, Aalborg University Hospital, Aalborg, Denmark
[3]Department of Mathematics, Aarhus University, Aarhus, Denmark

# 1  Supplementary Methods

## 1.1  Exposure Frobenius Distance

The difference between two matched weight matrices $\tilde{\boldsymbol{W}}, \hat{\boldsymbol{W}} \in \mathbb{R}_+^{K \times M}$ is measured as the average Frobenius distance between the two matrices:

$$L_F(\tilde{\boldsymbol{W}}, \hat{\boldsymbol{W}}) = \frac{1}{K \cdot M} ||\tilde{\boldsymbol{W}} - \hat{\boldsymbol{W}}||_F, \tag{S1}$$

where the weights are rearranged to match using the signature matching found in the corresponding signature sets.

## 1.2  Test error

The test error is calculated as the loss between the test data and its reconstruction. First, we divide the data into a training dataset $\boldsymbol{V}_{\text{train}}$ and a test dataset $\boldsymbol{V}_{\text{test}}$. The basis matrix $\hat{\boldsymbol{H}}$, and the weight matrix, $\hat{\boldsymbol{W}}_{\text{train}}$, of the training dataset are estimated. Then, for the test dataset the basis matrix from the training set, $\hat{\boldsymbol{H}}$, is fixed and the weight matrix is estimated using non-negative least squares [Lawson and Hanson, 1976]:

$$\hat{\boldsymbol{W}}_{\text{test}} = \arg \min_{\boldsymbol{W} \geq 0} ||\hat{\boldsymbol{H}}\boldsymbol{W} - \boldsymbol{V}_{\text{test}}||_F, \tag{S2}$$

from which the reconstructed test dataset is defined as $\hat{\boldsymbol{V}}_{\text{test}} := \hat{\boldsymbol{H}}\hat{\boldsymbol{W}}_{\text{test}}$. This representation is the test dataset reconstructed with the basis found in the training data. Now, the test error is calculated as

$$L_F(\boldsymbol{V}_{\text{test}}, \hat{\boldsymbol{V}}_{\text{test}}) = \frac{1}{M \cdot N} ||\boldsymbol{V}_{\text{test}} - \hat{\boldsymbol{V}}_{\text{test}}||_F. \tag{S3}$$

## 1.3  Choosing the number of basis vectors

Let $K$ denote the number of basis vectors used for extraction. The test error for each method and $K \in \{2, \ldots, \mathcal{K}\}$ is calculated using Algorithm S1 which is similar to the SigMos algorithm [Pelizzola, Laursen, and Hobolth, 2023]. The main difference being that Algorithm S1 samples with replacement and SigMos samples without replacement.

We consider the appropriate value of $K$ to be the value of $K$, where the test error does not change significantly when increasing $K$. To obtain this, iterative tests are performed on whether the test errors for a given $K$ can be assumed to have the same distribution as the test errors from the following $K + 1$. This iterative procedure continues until a test for similar distributions does not fail anymore and $K$ at

non-failure values is determined as the optimal number of basis vectors.

Since the same train/test splits are used to calculate the test errors for each $K$, a two-sample paired Wilcoxon's test is used to test whether the sample from a given $K$ can be assumed to be distributed as the sample from $K + 1$. A p-value of $p = 0.05$ is used as significance level. The exact procedure can be seen in Algorithm S2.

To ensure comparability in the quantitative analyses, we choose to use the same number of basis vectors for all methods. In cases where Algorithm S2 chooses different $K$'s for each method, the final number of basis vectors is determined the weighted average of each method rounded to the nearest integer:

$$K_{\text{all}} = \left\lceil \frac{K_{\text{NMF}} + \frac{1}{2} K_{\text{C-NMF}} + \frac{1}{2} K_{\text{AE-NMF}}}{2} \right\rceil, \tag{S4}$$

with AE-NMF and C-NMF weighted by $1/2$, since these are more likely to agree on the number of bases as they are variants of the same base model.

---

**Algorithm S1:** Test error as a function of number of basis vectors.

Input $\boldsymbol{V} \in \mathbb{R}_+^{M \times N}$ and $nsims$;
**for** $i \in \{1, \ldots, nsims\}$ **do**
    Sample $\text{idx}_{\text{train}}$ from $\{1, \ldots, N\}$ with replacement;
    $\boldsymbol{V}_{\text{train}} = \boldsymbol{V}[, \text{idx}_{\text{train}}] \in \mathbb{R}_+^{M \times N_{\text{train}}}$;
    $\boldsymbol{V}_{\text{test}} = \boldsymbol{V}[, -\text{idx}_{\text{train}}] \in \mathbb{R}_+^{M \times N_{\text{test}}}$;
    **for** $model \in \{NMF, C-NMF, AE-NMF\}$ *and* $K \in \{2, \ldots, \mathcal{K}\}$ **do**
        Fit $\hat{\boldsymbol{H}}, \hat{\boldsymbol{W}}_{\text{train}}$ using $K$-dimensional model on $\boldsymbol{V}_{\text{train}}$;
        Fit $\hat{\boldsymbol{W}}_{\text{train}} = \arg\min_{\boldsymbol{W} \geq 0} ||\hat{\boldsymbol{H}}\boldsymbol{W} - \boldsymbol{V}_{\text{test}}||_F^2$ ;
        Estimate $\hat{\boldsymbol{V}}_{\text{test}} = \hat{\boldsymbol{H}}\hat{\boldsymbol{W}}_{\text{test}}$;
        Calculate $x_{i,K,model} = L_F(\boldsymbol{V}_{\text{test}}, \hat{\boldsymbol{V}}_{\text{test}})/(M \cdot N_{\text{test}})$;
    **end**
**end**
**return** $\{x_{i,K,model}\}$ for $i \in \{1, \ldots, nsims\}$, $K \in \{2, \ldots, \mathcal{K}\}$ and
$model \in \{NMF, C-NMF, AE-NMF\}$

---

**Algorithm S2:** Iterative tests for the appropriate number of basis vectors.

Input test error sample $x_{1,K}, \ldots, x_{nsims,K}$ for $K = 2, \ldots, \mathcal{K}$;
**for** $K \in \{2, \ldots, \mathcal{K} - 1\}$ **do**
    Calculate $p$ as the p-value of Wilcoxon's two sample paired test of $\{x_{i,K}\}_{i=1,\ldots,nsims}$ and
    $\{x_{i,K+1}\}_{i=1,\ldots,nsims}$;
    **if** $p < p_{val}$ **then**
       | Return $K$
    **end**
    **else**
       | Continue
    **end**
**end**
Return $K$

---

## 1.4 Non-negativity in AE-NMF

There exist several ways of enforcing non-negativity in AE-NMF. An overview of these considered in this paper is provided in Table S1. Pei et al. (2020) suggested a ReLU activation function on the encoding layer and a softmax activation function on the decoding layer (ReLU + SM). Özer et al. (2022) found that using ReLU on the encoding layer while using a projected gradient on $\boldsymbol{W}_{\text{dec}}$ promoted convergence (ReLU + PG) Lin, 2007. Other methods enforce non-negativity with multiplicative updating steps. Khatib et al. (2018), e.g., suggested updating using exponential gradient descent and Zunner (2021) suggested

| Name | Formula | Source |
|---|---|---|
| ReLU + PG | $ReLU(\boldsymbol{V}\boldsymbol{W}_{\mathrm{enc}})\lvert\boldsymbol{W}_{\mathrm{dec}}\rvert$ | Lin, 2007 |
| ReLU + SM | $Softmax(ReLU(\boldsymbol{V}\boldsymbol{W}_{\mathrm{enc}}\boldsymbol{W}_{\mathrm{dec}}))$ | Pei et al., 2020 |
| EGD | $\boldsymbol{W}_{t+1} = \boldsymbol{W}_t \exp(-\eta\frac{\partial L_t}{\partial \boldsymbol{W}_t})$ | Khatib et al., 2018 |
| MU | $\boldsymbol{W}_{\mathrm{dec},nk}^{(t+1)} = \boldsymbol{W}_{\mathrm{dec},nk}^{(t)} \frac{(\boldsymbol{V}(\boldsymbol{W}_{\mathrm{enc}}\boldsymbol{V})^T)_{nk}}{(\boldsymbol{W}_{\mathrm{dec}}^{(t)}\boldsymbol{W}_{\mathrm{enc}}\boldsymbol{V}(\boldsymbol{W}_{\mathrm{enc}}\boldsymbol{V})^T)_{nk}}$ $\boldsymbol{W}_{\mathrm{enc},kn}^{(t+1)} = \boldsymbol{W}_{\mathrm{enc},kn}^{(t)} \frac{((\boldsymbol{W}_{\mathrm{dec}}^T\boldsymbol{V})\boldsymbol{V}^T)_{kn}}{((\boldsymbol{W}_{\mathrm{dec}}^T\boldsymbol{W}_{\mathrm{dec}}\boldsymbol{W}_{\mathrm{enc}}^{(t)}\boldsymbol{V})\boldsymbol{V}^T)_{kn}}$ | Zunner, 2021 |
| PG/FP PG | $\boldsymbol{V}ReLU(\boldsymbol{W}_{\mathrm{enc}})ReLU(\boldsymbol{W}_{\mathrm{dec}})$ | - |
| Abs/FP Abs | $\boldsymbol{V}\lvert\boldsymbol{W}_{\mathrm{enc}}\rvert\lvert\boldsymbol{W}_{\mathrm{dec}}\rvert$ | - |

Table S1: Overview of non-negative constraining methods for autoencoders.

multiplicative updates using a cleverly constructed learning rate. These methods are all excluded in this study, since methods enforcing non-negativity using activation functions on the network layers instead of the weight matrices (ReLU + PG and ReLU + SM) allow $\boldsymbol{W}_{\mathrm{enc}}$ and/or $\boldsymbol{W}_{\mathrm{dec}}$ to assume negative values, losing the equivalence with convex NMF. Methods based on multiplicative updates (EGD and MU) are also excluded in order to focus on the conventional additive updating scheme of autoencoders to contrast the multiplicative updating scheme of C-NMF and NMF

This study considers a projected gradient on both $\boldsymbol{W}_{\mathrm{enc}}$ and $\boldsymbol{W}_{\mathrm{dec}}$ either after each update (PG) or as a part of the forward pass (FP PG), and taking the absolute values of $\boldsymbol{W}_{\mathrm{enc}}$ and $\boldsymbol{W}_{\mathrm{dec}}$ after each update (Abs) or as a part of the forward pass (FP Abs).

## 2 Supplementary Results

### 2.1 Non-negativity in the autoencoder

The schema enforcing non-negativity in the encoding and decoding matrices, $\boldsymbol{W}_{\mathrm{enc}}$ and $\boldsymbol{W}_{\mathrm{dec}}$, in AE-NMF was settled by evaluating primarily the reconstruction error and secondarily the ACS with a corresponding signature set from C-NMF. These performance measures were chosen to obtain high reconstruction accuracy and similarity with C-NMF.
Dividing the 523 ovarian cancer patients' mutational profiles into 80/20 training/test set splits, 4 signatures were extracted from the training set using C-NMF and AE-NMF and each of the four non-negativity schemes (PG, FP PG, Abs, FP Abs). The training and test error for each non-negativity scheme in AE-NMF and the ACS between the C-NMF signatures and each of the AE-NMF signatures were calculated after extraction. This process was repeated for 50 train/test splits. A scatter plot of the ACS against the test errors as well as boxplots of the ACS and training and test errors for each non-negativity scheme across the 50 splits can be seen in Figure S2.
From these analyses it is observed that using a projected gradient (PG) or taking the absolute value in the forward pass (FP Abs) of the weight matrices perform similarly in both test error and ACS while outperforming the other methods. Methods enforcing non-negativity after each update (PG and Abs) had lower in-sample errors than methods enforcing non-negativity as a part of the forward pass (FP PG and FP Abs). Though PG appears to perform the best overall, it does experience problems by setting entire vectors to zero and being unable to exit this value again. Thus, all analyses will be carried out by taking the absolute value of negative weights in the forward pass (FP Abs) to enforce non-negativity in the weight matrices in AE-NMF.
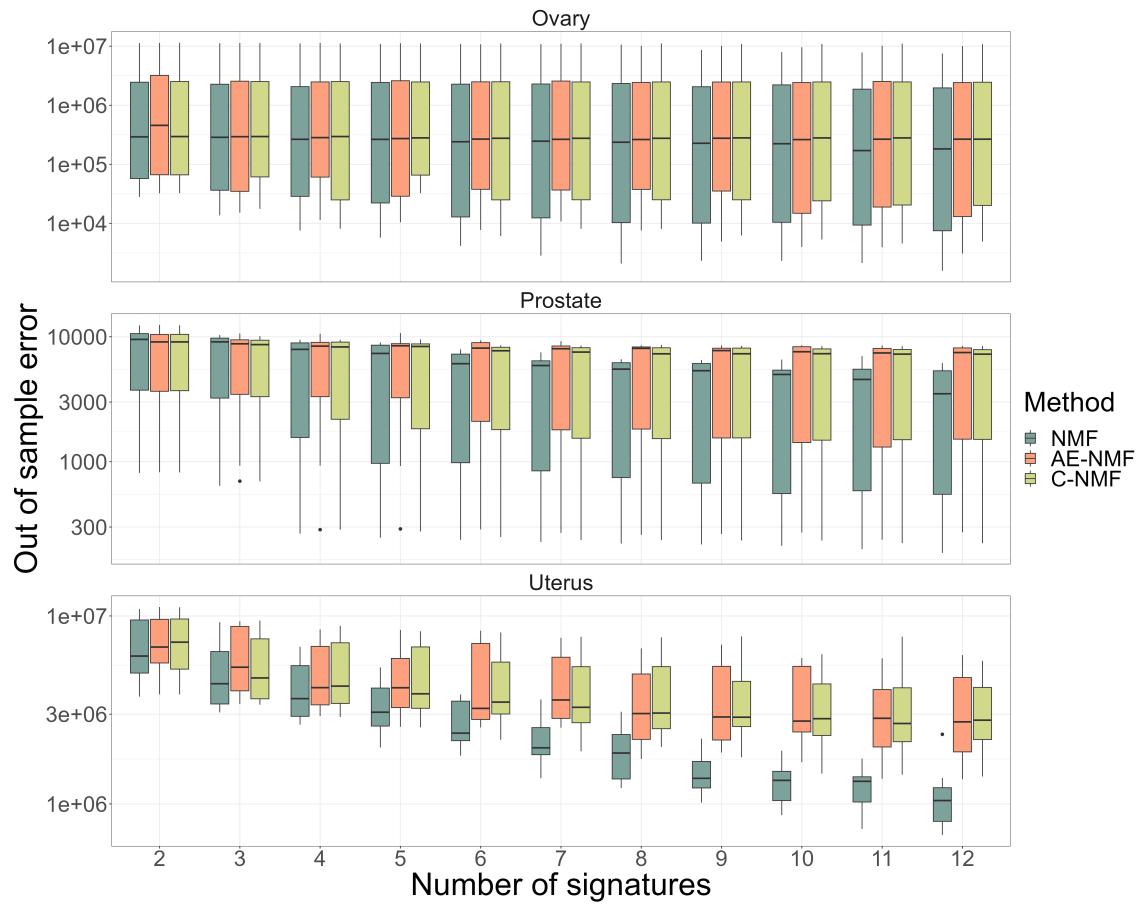
Figure S1: Boxplots of the test errors of the ten bootstrap train/test splits of the ovary, prostate and uterus cohorts against the number of signatures used in the extraction. Boxes are colored corresponding to the method used.
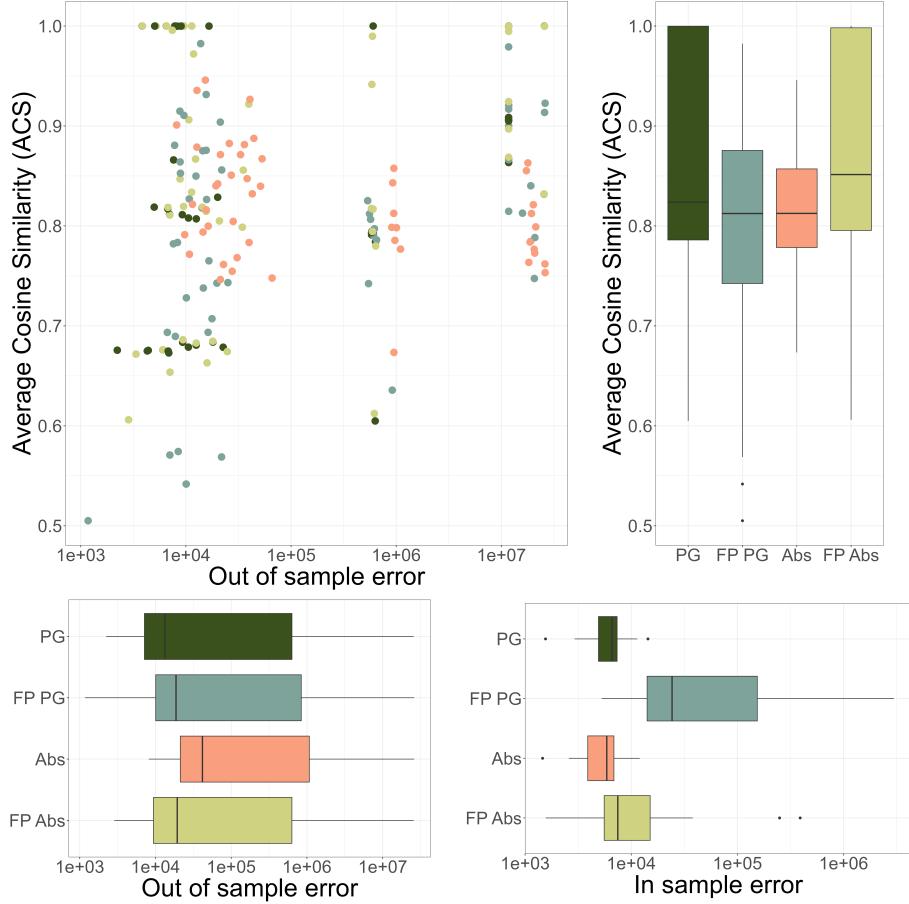
Figure S2: Top left: Scatter plot of the average cosine similarity (ACS) between C-NMF and AE-NMF for each non-negativity constraint against the out-of-sample error of AE-NMF for each non-negativity constraint. Top right: Boxplot of the ACS between C-NMF and AE-NMF for each non-negativity constraint. Bottom plots: Boxplots of the reconstruction errors for each non-negativity constraint of AE-NMF for the test sets (left) and the training sets (right).

Figure S3: Boxplots of the training- and test errors of 30 train/test splits of the ovary, prostate, and uterus cohorts. C-NMF and AE-NMF errors resulting from the same splits are connected by a line. The boxes are colored corresponding to the method used and the red diamond depicts the mean error. The y-axis is on a log 10 scale.
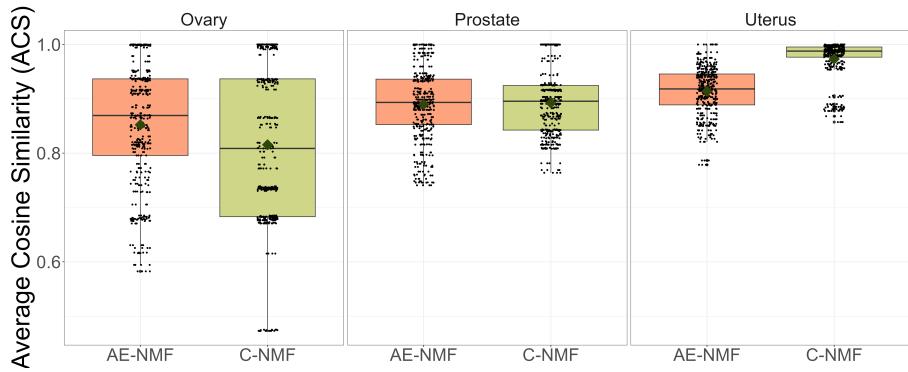


Figure S4: Box plots of the ACS between each combination of signatures extracted using a C-NMF or AE-NMF method across the 30 splits of data matrix for the ovary, prostate, and uterus cohort. Average consistency is marked by a red diamond.

| Model 1 | | C-NMF | C-NMF | AE-NMF |
|---------|---|-------|-------|--------|
| Model 2 | | AE-NMF | NMF | NMF |
| Data | K | $ACS(\hat{\boldsymbol{H}}_{\text{model 1}}, \hat{\boldsymbol{H}}_{\text{model 2}})$ | | |
| Ovary | 4 | **0.85** | 0.80 | 0.83 |
| Prostate | 5 | 0.90 | **0.95** | 0.83 |
| Uterus | 8 | **0.91** | 0.81 | 0.84 |
| Data | K | $L_F(\hat{\boldsymbol{W}}_{\text{model 1}}, \hat{\boldsymbol{W}}_{\text{model 2}})$ | | |
| Ovary | 4 | $\mathbf{1.1 \cdot 10^9}$ | $\mathbf{1.1 \cdot 10^9}$ | $1.2 \cdot 10^9$ |
| Prostate | 5 | $1.6 \cdot 10^7$ | $1.5 \cdot 10^7$ | $\mathbf{1.2 \cdot 10^7}$ |
| Uterus | 8 | $2.0 \cdot 10^{10}$ | $\mathbf{1.4 \cdot 10^{10}}$ | $1.5 \cdot 10^{10}$ |

Table S2: Average cosine similarity (ACS) and Frobenius distance when pairwisely comparing the signatures and exposures, respectively, found using C-NMF, AE-NMF, and NMF on the 30 splits of all training cohorts. Highest ACS and lowest Frobenius distance are highlighted in bold.

## 2.2 COSMIC Validation

Figures S5, S6, and S7 show the matched COSMIC signatures for all 30 train/test splits of the ovary, prostate, and uterus cohort, respectively. Each split is colored by the matched COSMIC signature, and the cosine similarity between the extracted signature and its COSMIC match is reported in each box.

### 2.2.1 Runtime

While NMF and C-NMF use multiplicative updating steps in their optimization AE-NMF uses additive, gradient descent-based updating steps. In Figure S8 the mean ($\pm$ SD) running time for signature extractions over the bootstrap samples is plotted against the number of signatures used in the extraction. We see that AE-NMF is consistently the fastest converging with the difference becoming more expressed as the number of signatures in the extraction increases. Optimizing with C-NMF becomes especially time consuming as the size of the factor matrices increases.

It it worth noting that AE-NMF updates are made with the Adam optimizer in the Python library PyTorch which is professionally constructed to be as time efficient as possible. In contrast, the code for the NMF and C-NMF updates used for this paper was written by the authors and without taking efficiency into special consideration. Therefore to be able to properly conclude on the running time of these algorithms more systematic comparisons are needed.
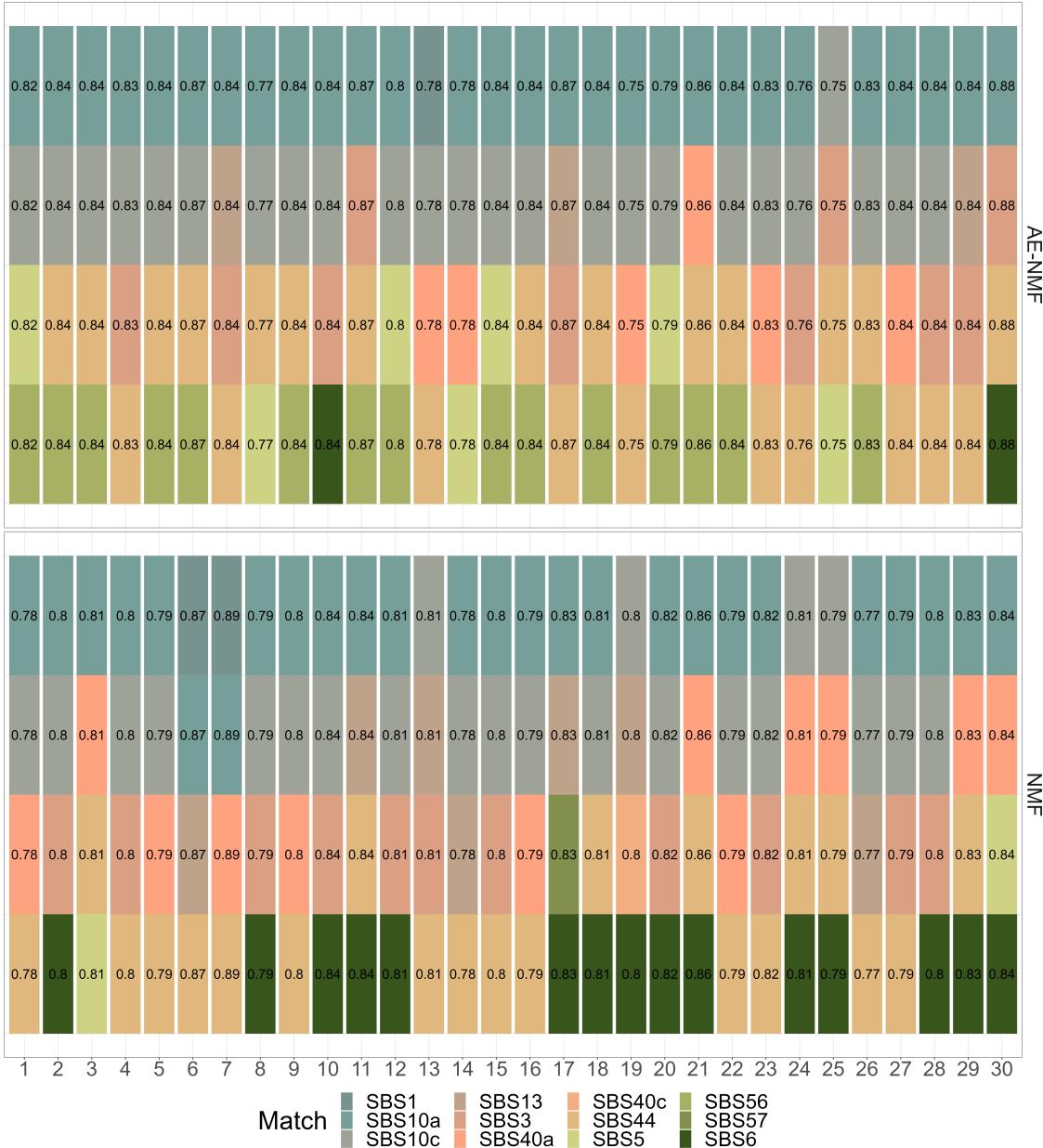
Figure S5: Cosine similarity between the four AE-NMF (top) or NMF (bottom) signatures extracted on the GEL ovary cohort and the matched COSMIC v.3.4 signature for each of the 30 train/test splits. The boxes are colored according to the matched COSMIC signature. The cosine similarity between each extracted signature and its COSMIC match is depicted within each box
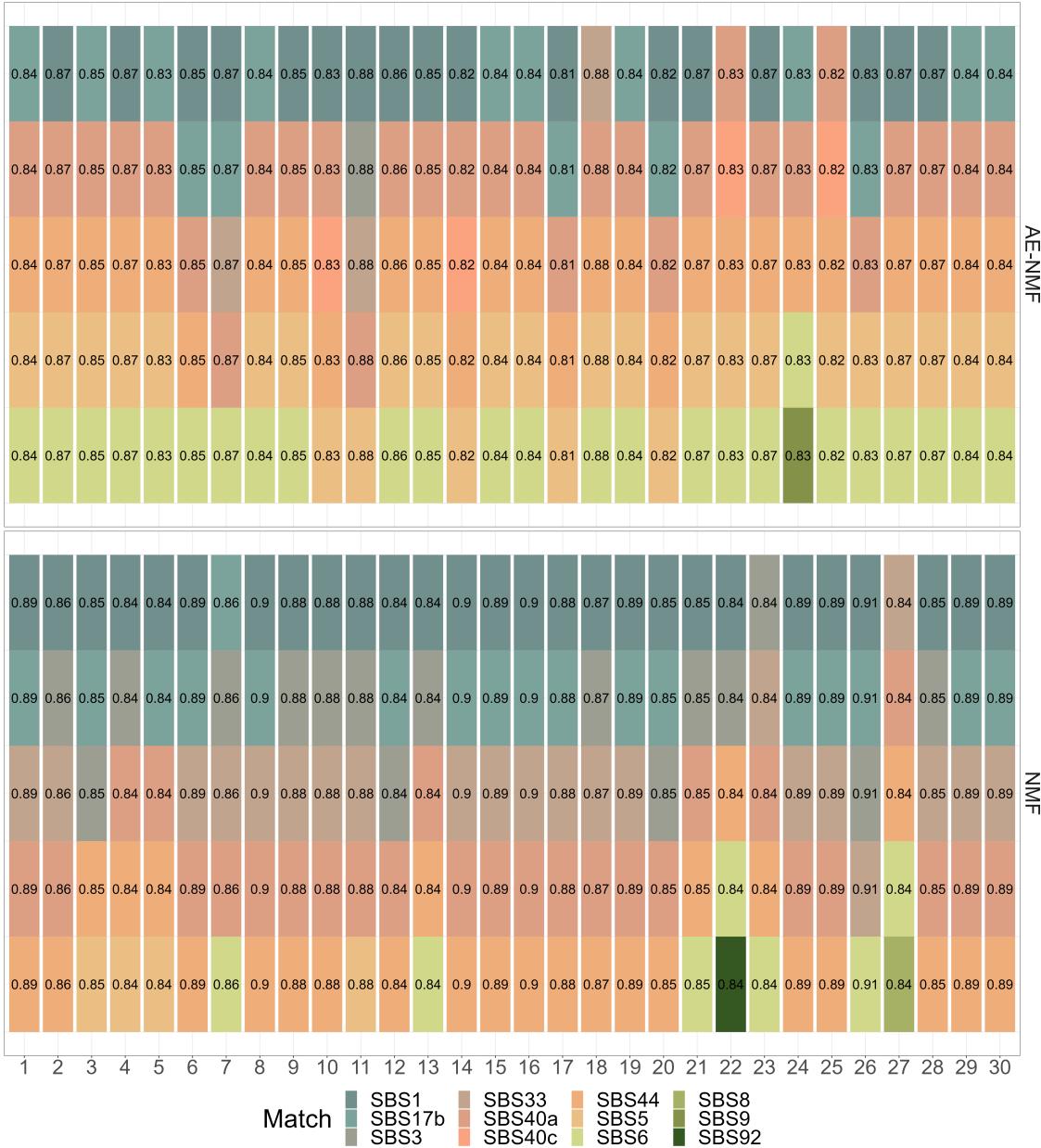
Figure S6: Cosine similarity between the 5 AE-NMF (top) or NMF (bottom) signatures extracted on the GEL ovary prostate and the matched COSMIC v.3.4 signature for each of the 30 train/test splits. The boxes are colored according to the matched COSMIC signature. The cosine similarity between each extracted signature and its COSMIC match is depicted within each box
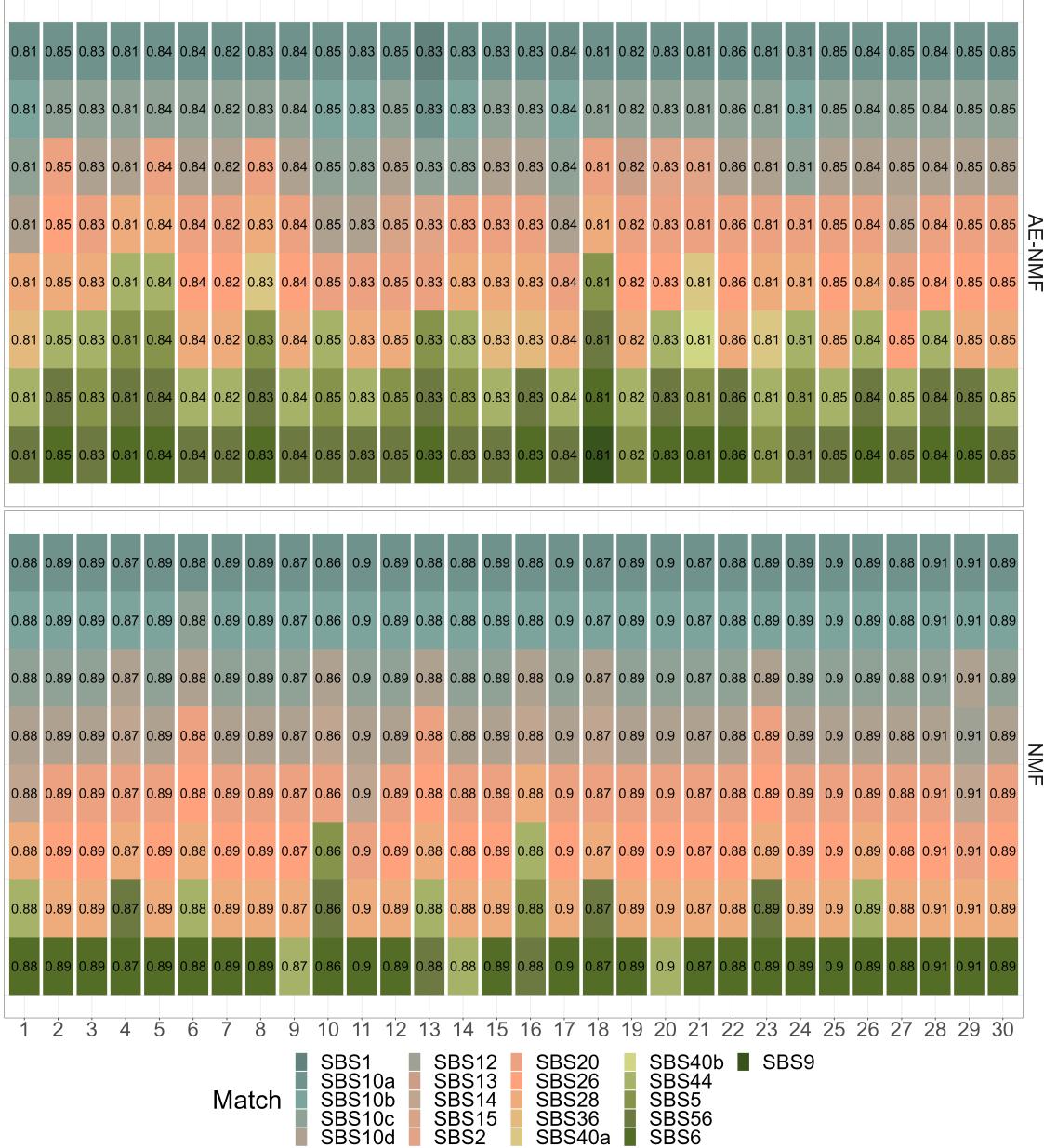
Figure S7: Cosine similarity between the 8 AE-NMF (top) or NMF (bottom) signatures extracted on the GEL uterus cohort and the matched COSMIC v.3.4 signature for each of the 30 train/test splits. The boxes are colored according to the matched COSMIC signature. The cosine similarity between each extracted signature and its COSMIC match is depicted within each box. The cosine similarity between each extracted signature and its COSMIC match is depicted within each box
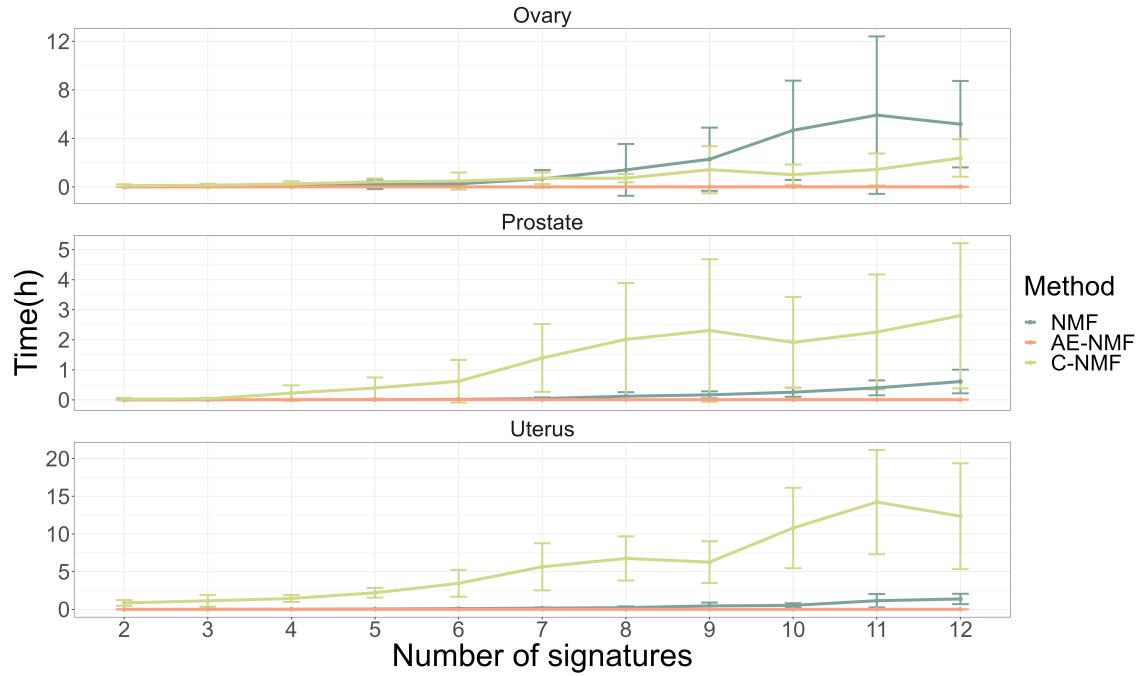
Figure S8: Mean (± SD) running time in hours over the bootstrap samples is plotted against the number of signatures ($K$) used for the extraction of the ovary, prostate, and uterus cohorts by each method in the bootstrap analyses to determine the number of signatures.

# References

Khatib, Alaa El et al. (2018). "Nonnegative Matrix Factorization Using Autoencoders And Exponentiated Gradient Descent". In: *Int. Joint Conf. on Neural Networks (IJCNN)*, pp. 1–8.

Lawson, Charles L. and Richard J. Hanson (1976). *Solving least squares problems*. SIAM.

Lin, Chih-Jen (2007). "Projected Gradient Methods for Nonnegative Matrix Factorization". In: *Neural Computation*, pp. 2756–2779.

Pei, Guangsheng et al. (2020). "Decoding whole-genome mutational signatures in 37 human pan-cancers by denoising sparse autoencoder neural network". In: *Oncogene* 39, 5031–5041.

Pelizzola, Marta, Ragnhild Laursen, and Asger Hobolth (2023). "Model selection and robust inference of mutational signatures using Negative Binomial non-negative matrix factorization". In: *BMC Bioinformatics* 24.

Zunner, Tim (2021). "Neural Networks with Nonnegativity Constraints for Decomposing Music Recordings". MA thesis. Friedrich-Alexander-Universität Erlangen-Nürnberg.

Özer, Yigitcan et al. (2022). "Investigating Nonnegative Autoencoders for Efficient Audio Decomposition". In: *30th European Signal Processing Conf. (EUSIPCO)*, pp. 254–258.