
Learning Binary Latent Variable Models: A Tensor Eigenpair Approach

Ariel Jaffe¹ Roi Weiss¹ Shai Carmi² Yuval Kluger³ Boaz Nadler¹

Abstract

Latent variable models with hidden binary units appear in various applications. Learning such models, in particular in the presence of noise, is a challenging computational problem. In this paper we propose a novel spectral approach to this problem, based on the eigenvectors of both the second order moment matrix and third order moment tensor of the observed data. We prove that under mild non-degeneracy conditions, our method consistently estimates the model parameters at the optimal parametric rate. Our tensor-based method generalizes previous orthogonal tensor decomposition approaches, where the hidden units were assumed to be either statistically independent or mutually exclusive. We illustrate the consistency of our method on simulated data and demonstrate its usefulness in learning a common model for population mixtures in genetics.

1. Introduction

In this paper we propose a spectral method for learning the following binary latent variable model, shown in Figure 1. The hidden layer, $\mathbf{h} = (h_1, \dots, h_d)$, consists of d binary random variables with an unknown joint distribution $P_{\mathbf{h}} : \{0, 1\}^d \rightarrow [0, 1]$. The observed vector $\mathbf{x} \in \mathbb{R}^m$ with $m \geq d$ features is modeled as

$$\mathbf{x} = W^T \mathbf{h} + \sigma \boldsymbol{\xi}, \quad (1)$$

where $W \in \mathbb{R}^{d \times m}$ is an unknown weight matrix assumed to be full rank d . Here, $\sigma \geq 0$ is the noise level and $\boldsymbol{\xi}$ is an additive noise vector independent of \mathbf{h} , whose m coordinates are all i.i.d. zero mean and unit variance random variables.

¹Dept. of Computer Science and Applied Mathematics, Weizmann Institute of Science, Rehovot 7610001, Israel. ²Braun School of Public Health and Community Medicine, The Hebrew University of Jerusalem, Jerusalem 9112102, Israel. ³Program of Applied Mathematics, Yale University, New Haven, CT 06511, USA. Correspondence to: Ariel Jaffe <ariel.jaffe@weizmann.ac.il>, Roi Weiss <roi.weiss@weizmann.ac.il>.

For simplicity we assume it is Gaussian, though our method can be modified to handle other noise distributions.

The model in (1) appears, for example, in overlapping clustering (Banerjee et al., 2005; Baadel et al., 2016), in various problems in bioinformatics (Segal et al., 2002; Becker et al., 2011; Slawski et al., 2013), and in blind source separation (Van der Veen, 1997). A special instance of model (1) is the Gaussian-Bernoulli restricted Boltzmann machine (G-RBM) where the distribution $P_{\mathbf{h}}$ is further assumed to have a parametric energy-based structure (Hinton & Salakhutdinov, 2006; Cho et al., 2011; Wang et al., 2012). G-RBMs were used, e.g., in modeling human motion (Taylor et al., 2007) and natural image patches (Melchior et al., 2017).

Given n i.i.d. samples $\mathbf{x}_1, \dots, \mathbf{x}_n$ from model (1), the goal is to estimate the weight matrix W . A common approach for learning W is by maximum likelihood. As this function is non-convex, common optimization schemes include the EM algorithm and alternating least squares (ALS). In addition, several works developed iterative methods specialized to G-RBMs (Hinton, 2010; Cho et al., 2011). All these methods, however, often lack consistency guarantees and may not be well suited for large datasets due to their potential slow convergence. This is not surprising, as learning W under model (1) is believed to be computationally hard; see for example Mossel & Roch (2005).

Over the past years, several works considered variants and specific instances of model (1) under additional assumptions on the distribution $P_{\mathbf{h}}$ or on the weight matrix W . For example, when $P_{\mathbf{h}}$ is a product distribution, the learning problem becomes that of independent component analysis (ICA) with binary signals (Hyvärinen et al., 2004). In this case, several methods were derived for estimating W and under suitable non-degeneracy conditions were proven to be both computationally efficient and statistically consistent (Shalvi & Weinstein, 1993; Frieze et al., 1996; Regalia & Kofidis, 2003; Hyvärinen et al., 2004; Anandkumar et al., 2014; Jain & Oh, 2014). Similarly, when the hidden units are mutually exclusive, namely $P_{\mathbf{h}}$ has support $\mathbf{h} \in \{\mathbf{e}_i\}_{i=1}^d$, the model is a Gaussian mixture (GMM) with d spherical components with linearly independent means. Efficient and consistent algorithms were derived for this case as well (Moitra & Valiant, 2010; Anandkumar et al., 2012a;b; Hsu & Kakade, 2013). Among those, most relevant to this work

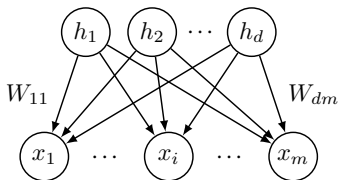


Figure 1. The binary latent variable model.

are orthogonal tensor decomposition methods (Anandkumar et al., 2014). Interestingly, these methods can learn some additional latent models, with hidden units that are not necessarily binary, such as Dirichlet allocation and other correlated topic models (Arabshahi & Anandkumar, 2017).

Learning W given the observed data $\{x_j\}_{j=1}^n$ can also be viewed as a noisy *matrix factorization* problem. If W is known to be non-negative, then various non-negative matrix factorization methods can be used. Moreover, under appropriate conditions, some of these methods were proven to be computationally efficient and consistent (Donoho & Stodden, 2004; Arora et al., 2012). For general full rank W , the matrix factorization method in Slawski et al. (2013) (SHL) exactly recovers W when $\sigma = 0$ with a runtime exponential in d . This method, however, can handle only low levels of noise and has no consistency guarantees when $\sigma > 0$.

A tensor eigenpair approach In this paper we propose a novel spectral method for learning W which is based on the eigenvectors of both the second order moment matrix and the third order moment tensor of the observed data. We prove that our method is consistent under mild non-degeneracy conditions and achieves the parametric rate $O_P(n^{-\frac{1}{2}})$ for any noise level $\sigma \geq 0$.

The non-degeneracy conditions we pose are significantly weaker than those required by the previous tensor decomposition methods mentioned above. In particular, their assumptions and resulting methods can be viewed as specific cases of our more general approach.

Similarly to the matrix factorization method in Slawski et al. (2013), our algorithm has runtime linear in n , polynomial in m , and in general exponential in d . With our current Matlab implementation, most of the runtime is spent on computing the eigenpairs of a $d \times d \times d$ tensor. Practically, our method, implemented without any particular optimization, can learn a model with 12 hidden units in less than ten minutes on a standard PC. Furthermore, the overall runtime can be significantly reduced, since the step of computing the tensor eigenpairs can be embarrassingly parallelized.

Paper outline In the next section we give necessary background on tensor eigenpairs. In Section 3 we introduce our

method in the case $\sigma = 0$. The case $\sigma \geq 0$ is treated in Section 4. Experiments with our method and comparison to other approaches appear in Section 5. All proofs are deferred to the supplementary material.

2. Preliminaries

Notation Denote $[d] = \{1, \dots, d\}$ and e_i as the i -th unit vector. We slightly abuse notation and view a matrix W also as the set of its columns, namely $w \in W$ is some column of W and $\text{span}(W)$ is the span of all its columns. The unit sphere is denoted by $\mathbb{S}_{d-1} = \{u \in \mathbb{R}^d : \|u\| = 1\}$.

A tensor $\mathcal{T} \in \mathbb{R}^{d \times d \times d}$ is symmetric if $\mathcal{T}_{ijk} = \mathcal{T}_{\pi(i,j,k)}$ for all permutations π of i, j, k . Here, we consider only symmetric tensors. \mathcal{T} can also be seen as a *multi-linear* operator: for matrices W^1, W^2, W^3 with $W^i \in \mathbb{R}^{d \times d_i}$, the tensor-mode product, denoted $\mathcal{T}(W^1, W^2, W^3)$, is a $d_1 \times d_2 \times d_3$ tensor whose (i_1, i_2, i_3) -th entry is

$$\sum_{j_1, j_2, j_3 \in [d]} W_{j_1 i_1}^1 W_{j_2 i_2}^2 W_{j_3 i_3}^3 \mathcal{T}_{j_1 j_2 j_3}.$$

Tensor eigenpairs Several types of eigenpairs of a tensor have been proposed. Here, we consider the following definition, termed Z -eigenpairs by Qi (2005) and l_2 -eigenpairs by Lim (2005). Henceforth we just call them eigenpairs.

Definition 1. $(u, \lambda) \in \mathbb{R}^d \times \mathbb{R}$ is an eigenpair of \mathcal{T} if

$$\mathcal{T}(I, u, u) = \lambda u \quad \text{and} \quad \|u\| = 1. \quad (2)$$

Note that if (u, λ) is an eigenpair then the eigenvalue is simply $\lambda = \mathcal{T}(u, u, u)$. In addition, $(-u, -\lambda)$ is also an eigenpair. Following common practice we treat these two pairs as one and make the convention that $\lambda \geq 0$.

In contrast to the matrix case, the number of eigenvalues $\{\lambda\}$ of a tensor $\mathcal{T} \in \mathbb{R}^{d \times d \times d}$ can be much larger than d . As shown by Cartwright & Sturmfels (2013), for a $d \times d \times d$ tensor, there can be at most $2^d - 1$ of them. With precise definitions appearing in Cartwright & Sturmfels (2013), for a *generic* tensor, all its eigenvalues have multiplicity one and the number of eigenpairs $\{(u, \lambda)\}$ is at most $2^d - 1$.

In principle, enumerating the set of all eigenpairs of a general symmetric tensor is a #P problem (Hillar & Lim, 2013). Nevertheless, several methods have been proposed for computing at least *some* eigenpairs, including iterative higher-order power methods (Kolda & Mayo, 2011; 2014), homotopy continuation (Chen et al., 2016), semidefinite programming (Cui et al., 2014), and iterative Newton-based methods (Jaffe et al., 2017; Guo et al., 2017). We conclude this section with the definition of *Newton-stable* eigenpairs (Jaffe et al., 2017) which are most relevant to our work.

Newton-stable eigenpairs Equivalently to (2), eigenpairs of \mathcal{T} can also be characterized by the function $g : \mathbb{R}^d \rightarrow \mathbb{R}^d$,

$$g(\mathbf{u}) = \mathcal{T}(I, \mathbf{u}, \mathbf{u}) - \mathcal{T}(\mathbf{u}, \mathbf{u}, \mathbf{u}) \cdot \mathbf{u}. \quad (3)$$

It is easy to verify that a pair (\mathbf{u}, λ) with $\|\mathbf{u}\| = 1$ is an eigenpair of \mathcal{T} if and only if $g(\mathbf{u}) = \mathbf{0}$ and $\lambda = \mathcal{T}(\mathbf{u}, \mathbf{u}, \mathbf{u})$. The stability of an eigenpair is determined by its Jacobian matrix $\nabla g(\mathbf{u}) \in \mathbb{R}^{d \times d}$, more precisely, by its projection into the $d - 1$ dimensional subspace orthogonal to \mathbf{u} . Formally, let $L_{\mathbf{u}} \in \mathbb{R}^{d \times (d-1)}$ be a matrix with $d - 1$ orthonormal columns that span the subspace orthogonal to \mathbf{u} and define the $(d - 1) \times (d - 1)$ projected Jacobian matrix

$$J_p(\mathbf{u}) = L_{\mathbf{u}}^\top \nabla g(\mathbf{u}) L_{\mathbf{u}}. \quad (4)$$

Definition 2. An eigenpair (\mathbf{u}, λ) of $\mathcal{T} \in \mathbb{R}^{d \times d \times d}$ is Newton-stable if the matrix $J_p(\mathbf{u})$ has full rank $d - 1$.

The homotopy continuation method in Chen et al. (2016) is guaranteed to compute all the Newton-stable eigenpairs of a tensor. Alternatively, all the Newton-stable eigenpairs can be computed by the iterative orthogonal Newton correction method (O-NCM) in Jaffe et al. (2017) as these are the attracting fixed points for this algorithm. Moreover, O-NCM converges to any Newton-stable eigenpair at a quadratic rate given a sufficiently close initial guess. Finally, for a generic tensor, all its eigenpairs are Newton-stable.

3. Learning in the noiseless case

To motivate our approach for estimating the matrix W it is instructive to first consider the ideal noiseless case where $\sigma = 0$. In this case, model (1) takes the form $\mathbf{x} = W^\top \mathbf{h}$. Our problem then becomes that of factorizing the observed matrix $X = [\mathbf{x}_1, \dots, \mathbf{x}_n] \in \mathbb{R}^{m \times n}$ of n samples into a product of real and binary low-rank matrices,¹

$$\text{Find } W \in \mathbb{R}^{d \times m}, H \in \{0, 1\}^{d \times n} \text{ s.t. } X = W^\top H. \quad (5)$$

To be able to recover W we first need conditions under which the decomposition of X into W and H is unique. Clearly, such a factorization can be unique at most up to a permutation of its components; we henceforth ignore this degeneracy. A sufficient condition for uniqueness, similar to the one posed in Slawski et al. (2013), is that H is *rigid*. Formally, $H \in \{0, 1\}^{d \times n}$ is rigid if any non-trivial linear combination of its rows yields a non-binary vector: $\forall \mathbf{u} \neq \mathbf{0}$,

$$\mathbf{u}^\top H \in \{0, 1\}^n \Leftrightarrow \mathbf{u} \in \{\mathbf{e}_i\}_{i=1}^d. \quad (6)$$

Condition (6) is satisfied, for example, when the columns of H include \mathbf{e}_i and $\mathbf{e}_i + \mathbf{e}_j$ for all $i \neq j \in [d]$. If there

¹Note that this is different from the problem known as ‘‘Boolean matrix factorization’’, where X and W are assumed to be binary as well; see Miettinen & Vreeken (2014) and references therein.

exists a positive constant $p_0 > 0$ such that $P_{\mathbf{h}}(\mathbf{e}_i) \geq p_0$ and $P_{\mathbf{h}}(\mathbf{e}_i + \mathbf{e}_j) \geq p_0$, then for a sample size $n > 2 \log(d)/p_0$ the matrix H is rigid with high probability.

The following proposition, similar in nature to the (affine constrained) uniqueness guarantee in Slawski et al. (2013), shows that under condition (6) the factorization in (5) is unique and fully characterized by the binary constraints.

Proposition 1. Let $X = W^\top H$ with $H \in \{0, 1\}^{d \times n}$ rigid and $W \in \mathbb{R}^{d \times m}$ full rank with $m \geq d$. Let $W^\dagger \in \mathbb{R}^{m \times d}$ be the unique right pseudo-inverse of W so $W W^\dagger = I_d$. Then W and H are unique and for all $\mathbf{v} \in \text{span}(X) \setminus \{\mathbf{0}\}$,

$$\mathbf{v}^\top X \in \{0, 1\}^n \Leftrightarrow \mathbf{v} \in W^\dagger. \quad (7)$$

Hence, under the rigidity condition (6), the matrix factorization problem in (5) is equivalent to the problem of finding the *unique* set $W^\dagger = \{\mathbf{v}_1^*, \dots, \mathbf{v}_d^*\} \subseteq \text{span}(X)$ of d non-zero vectors that satisfy the binary constraints $\mathbf{v}_i^{*\top} X \in \{0, 1\}^n$. The weight matrix is then $W = (W^\dagger)^\dagger$.

Algorithm outline We recover W^\dagger via a two step procedure. First, a finite set $V = \{\mathbf{v}_1, \mathbf{v}_2, \dots\} \subseteq \text{span}(X)$ of *candidate* vectors is computed with a guarantee that $W^\dagger \subseteq V$. Specifically, V is computed from the set of eigenpairs of a $d \times d \times d$ tensor, constructed from the low order moments of X . Typically, the size of V will be much larger than d , so in the second step V is *filtered* by selecting all $\mathbf{v} \in V$ that satisfy $\mathbf{v}^\top X \in \{0, 1\}^n$.

Before describing the two steps in more detail we first state the additional non-degeneracy conditions we pose. To this end, denote the unknown first, second, and third order moments of the latent binary vector \mathbf{h} by

$$\begin{aligned} \mathbf{p} &= \mathbb{E}[\mathbf{h}] \in \mathbb{R}^d, & C &= \mathbb{E}[\mathbf{h} \otimes \mathbf{h}] \in \mathbb{R}^{d \times d}, \\ C &= \mathbb{E}[\mathbf{h} \otimes \mathbf{h} \otimes \mathbf{h}] \in \mathbb{R}^{d \times d \times d}. \end{aligned} \quad (8)$$

Non-degeneracy conditions We assume the following:

- (I) H is rigid.
- (II) $\text{rank}(2C(I, I, \mathbf{e}_i) - C) = d$ for all $i \in [d]$.

Condition (I) implies that both $\text{rank}(H H^\top) = d$ and $\text{rank}(C) = d$. This in turn implies $p_i = \mathbb{E}[h_i] > 0$ for all $i \in [d]$ and that at most one variable h_i has $p_i = 1$. Such an ‘‘always on’’ variable can model a fixed bias to \mathbf{x} . As far as we know, condition (II) is new and its nature will become clear shortly.

We now describe each step of our algorithm in more detail.

Computing the candidate set To compute a set V that is guaranteed to include the columns of W^\dagger we make use of

the second and third order moments of \mathbf{x} ,

$$\begin{aligned} M &= \mathbb{E}[\mathbf{x} \otimes \mathbf{x}] \in \mathbb{R}^{m \times m}, \\ \mathcal{M} &= \mathbb{E}[\mathbf{x} \otimes \mathbf{x} \otimes \mathbf{x}] \in \mathbb{R}^{m \times m \times m}. \end{aligned} \quad (9)$$

Given a large number of samples $n \gg 1$, these can be easily and accurately estimated from the sample X . For simplicity, in this section we consider the population setting where $n \rightarrow \infty$, so M and \mathcal{M} are known exactly. M and \mathcal{M} are related to the unknown second and third order moments of \mathbf{h} in (8) via (Anandkumar et al., 2014)

$$M = W^\top C W, \quad \mathcal{M} = \mathcal{C}(W, W, W). \quad (10)$$

Since both C and W are full rank, the number of latent units can be deduced by $\text{rank}(M) = d$. Since C is positive definite, there is a whitening matrix $K \in \mathbb{R}^{m \times d}$ such that

$$K^\top M K = I_d. \quad (11)$$

Such a K can be computed, for example, by an eigen-decomposition of M . Although K is not unique, any $K \subseteq \text{span}(M)$ that satisfies (11) suffices for our purpose. Define the $d \times d \times d$ lower dimensional whitened tensor

$$\mathcal{W} = \mathcal{M}(K, K, K). \quad (12)$$

Denote the set of *eigenpairs* of \mathcal{W} by

$$U = \{(\mathbf{u}, \lambda) \in \mathbb{S}_{d-1} \times \mathbb{R}_+ : \mathcal{W}(\mathbf{u}, \mathbf{u}) = \lambda \mathbf{u}\}. \quad (13)$$

Our set of candidates is then

$$V = \{K\mathbf{u}/\lambda : (\mathbf{u}, \lambda) \in U \text{ with } \lambda \geq 1\} \subseteq \mathbb{R}^m. \quad (14)$$

The following lemma shows that under condition (I) the set V is guaranteed to contain the d columns of W^\dagger .

Lemma 1. *Let \mathcal{W} be the tensor in (12) corresponding to model (1) with $\sigma = 0$ and let V be as in (14). If condition (I) holds then $W^\dagger \subseteq V$. In particular, each $(\mathbf{u}_i, \lambda_i)$ in the set of d relevant eigenpairs*

$$U^* = \{(\mathbf{u}, \lambda) \in U : K\mathbf{u}/\lambda \in W^\dagger\} \quad (15)$$

has the eigenvalue $\lambda_i = 1/\sqrt{p_i} \geq 1$ where $p_i = \mathbb{E}[h_i] > 0$.

Computing the tensor eigenpairs By Lemma 1, we may construct a candidate set V that contains W^\dagger by first calculating the set U of eigenpairs of \mathcal{W} . Unfortunately, computing the set of all eigenpairs of a general symmetric tensor is computationally hard (Hillar & Lim, 2013). Moreover, besides the d columns of W^\dagger , the set V in (14) may contain many spurious candidates, as the number of eigenpairs of \mathcal{W} is typically $O(2^d) \gg d$ (Cartwright & Sturmfels, 2013).

Nevertheless, as discussed in Section 2, several methods have been proposed for computing *some* eigenpairs of a tensor under appropriate stability conditions. The following lemma highlights the importance of condition (II) for the stability of the eigenpairs in U^* . Note that conditions (I)-(II) do not depend on W , but only on the distribution of \mathbf{h} .

Lemma 2. *Let \mathcal{W} be the whitened tensor in (12) corresponding to model (1) with $\sigma = 0$. If conditions (I)-(II) hold, then all $(\mathbf{u}, \lambda) \in U^*$ are Newton-stable eigenpairs of \mathcal{W} .*

Hence, under conditions (I)-(II), the homotopy method in Chen et al. (2016), or alternatively the O-NCM with a sufficiently large number of random initializations (Jaffe et al., 2017), are guaranteed to compute a candidate set V which includes all the columns of W^\dagger . The next step is to extract W^\dagger out of V .

Filtering As suggested by Eq. (7) we select the subset of vectors $\bar{V} \subseteq V$ that satisfy the binary constraints,

$$\bar{V} = \{\mathbf{v} \in V : \mathbf{v}^\top X \in \{0, 1\}^n\}. \quad (16)$$

Indeed, under condition (I), Proposition 1 implies that $\bar{V} = W^\dagger$ and the weight matrix is thus $W = \bar{V}^\dagger$.

Algorithm 2 in Appendix C summarizes our method for the noiseless case and has the following recovery guarantee.

Theorem 1. *Let X be a matrix of n samples from model (1) with $\sigma = 0$. If conditions (I)-(II) hold, then the above method recovers W exactly.*

We note that when $\sigma = 0$ and conditions (I)-(II) hold for the *empirical* latent moments \hat{C} and $\hat{\mathcal{C}}$ (rather than C and \mathcal{C}), the above procedure *exactly* recovers W when M and \mathcal{M} are replaced by their finite sample estimates. The matrix factorization method SHL in Slawski et al. (2013) also exactly recovers W in the case $\sigma = 0$. While its runtime is also exponential in d , practically it may be much faster than our proposed tensor based approach. This is because SHL constructs a candidate set of size 2^d that can be computed by a suitable linear transformation of the *fixed* set $\{0, 1\}^d$, as opposed to our candidate set which is constructed by eigenpairs of a $d \times d \times d$ tensor. However, SHL does not take advantage of the large number of samples n , since only $m \times d$ sub-matrices of the $m \times n$ sample matrix X are used for constructing its candidate set. Indeed, in the noisy case where $\sigma > 0$, SHL has no consistency guarantees and as demonstrated by the simulation results in Section 5 it may fail at high levels of noise. In the next section we derive a modified version of our method that consistently estimates W for any noise level $\sigma \geq 0$.

4. Learning in the presence of noise

The method in Section 3 to estimate W is clearly inadequate when $\sigma > 0$. However, we now show that by making several adjustments, the two steps of computing the candidate set and its filtering can be both made robust to noise, yielding a consistent estimator of W for any $\sigma \geq 0$.

Computing the candidate set As in the case $\sigma = 0$, our goal in the first step is to compute a finite candidate set

$V_\sigma \subseteq \mathbb{R}^m$ that is guaranteed to contain accurate estimates for the d columns of W^\dagger . To this end, in addition to the second and third order moments M and \mathcal{M} in (9), we also consider the first order moment $\boldsymbol{\mu} = \mathbb{E}[\boldsymbol{x}]$ and define the following noise corrected moments,

$$\begin{aligned} M_\sigma &= M - \sigma^2 I_m, \\ \mathcal{M}_\sigma &= \mathcal{M} - \sigma^2 \sum_{i=1}^m \left(\boldsymbol{\mu} \otimes \mathbf{e}_i \otimes \mathbf{e}_i \right. \\ &\quad \left. + \mathbf{e}_i \otimes \boldsymbol{\mu} \otimes \mathbf{e}_i + \mathbf{e}_i \otimes \mathbf{e}_i \otimes \boldsymbol{\mu} \right). \end{aligned} \quad (17)$$

By assumption, the noise satisfies $\mathbb{E}[\xi_i^3] = 0$. Thus, similarly to the moment equations in (10), the modified moments in (17) are related to these of \boldsymbol{h} by (Anandkumar et al., 2014)

$$M_\sigma = W^\top C W, \quad \mathcal{M}_\sigma = \mathcal{C}(W, W, W). \quad (18)$$

Hence, if M_σ and \mathcal{M}_σ were known exactly, a candidate set V_σ that contains W^\dagger could be obtained exactly as in the noiseless case, but with M and \mathcal{M} replaced with M_σ and \mathcal{M}_σ ; namely, first calculate the whitening matrix K_σ such that $K_\sigma^\top M_\sigma K_\sigma = I_d$ and then compute the eigenpairs of the population whitened tensor

$$\mathcal{W}_\sigma = \mathcal{M}_\sigma(K_\sigma, K_\sigma, K_\sigma). \quad (19)$$

In practice, σ^2 , d , $\boldsymbol{\mu}$, M and \mathcal{M} are all unknown and need to be estimated from the sample matrix X . Assuming $m > d$, the parameters σ^2 and d can be consistently estimated, for example, by the methods in Kritchman & Nadler (2009). For simplicity, we assume they are known exactly. Similarly, $\boldsymbol{\mu}$, M , \mathcal{M} are consistently estimated by their empirical means, $\hat{\boldsymbol{\mu}}$, \hat{M} , and $\hat{\mathcal{M}}$. So, after computing the plugin estimates \hat{K}_σ such that $\hat{K}_\sigma^\top \hat{M}_\sigma \hat{K}_\sigma = I_d$ and $\hat{\mathcal{W}}_\sigma = \hat{\mathcal{M}}_\sigma(\hat{K}_\sigma, \hat{K}_\sigma, \hat{K}_\sigma)$, we compute the set \hat{U}_σ of eigenpairs of $\hat{\mathcal{W}}_\sigma$ and for some small $0 < \tau = O(n^{-\frac{1}{2}})$ take our candidate set as

$$\hat{V}_\sigma = \{\hat{K}_\sigma \boldsymbol{u} / \lambda : (\boldsymbol{u}, \lambda) \in \hat{U}_\sigma \text{ with } \lambda \geq 1 - \tau\}. \quad (20)$$

The following lemma shows that under conditions (I)-(II) the above procedure is stable to small perturbations. Namely, for perturbations of order $\delta \ll 1$ in \mathcal{W}_σ and K_σ , the method computes a candidate set \hat{V}_σ that contains a subset of d vectors that are $O(\delta)$ close to the columns of W^\dagger . Furthermore, these d vectors all correspond to Newton-stable eigenpairs of the perturbed tensor and are $\Omega(1)$ separated from the other candidates in \hat{V}_σ .

Lemma 3. *Let $K_\sigma, \mathcal{W}_\sigma$ be the population quantities in (19) and let $\hat{K}_\sigma, \hat{\mathcal{W}}_\sigma$ be their perturbed versions, inducing the candidate set \hat{V}_σ in (20). If conditions (I)-(II) hold, then there are $c, \delta_0, \delta_1 > 0$ such that for all $0 \leq \delta \leq \delta_0$ the following holds: If the perturbed versions satisfy*

$$\max\{\|\hat{\mathcal{W}}_\sigma - \mathcal{W}_\sigma\|_F, \|\hat{K}_\sigma - K_\sigma\|_F\} \leq \delta, \quad (21)$$

then any $\boldsymbol{v}^* \in W^\dagger$ has a unique $\hat{\boldsymbol{v}} \in \hat{V}_\sigma$ such that

$$\|\hat{\boldsymbol{v}} - \boldsymbol{v}^*\| \leq c\delta. \quad (22)$$

Moreover, $\hat{\boldsymbol{v}}$ corresponds to a Newton-stable eigenpair of $\hat{\mathcal{W}}_\sigma$ with eigenvalue $\lambda \geq 1 - c\delta$ and for all $\tilde{\boldsymbol{v}} \in \hat{V}_\sigma \setminus \{\hat{\boldsymbol{v}}\}$,

$$\|\tilde{\boldsymbol{v}} - \boldsymbol{v}^*\| \geq \delta_1 > 2c\delta. \quad (23)$$

The proof is based on the implicit function theorem (Hubbard & Hubbard, 2015); small perturbations to a tensor result in small perturbations to its Newton-stable eigenpairs.

Now, by the delta method, the plugin estimates \hat{K}_σ and $\hat{\mathcal{W}}_\sigma$ are both $O_P(n^{-\frac{1}{2}})$ close to their population quantities,

$$\begin{aligned} \|\hat{K}_\sigma - K_\sigma\|_F &= O_P(n^{-\frac{1}{2}}), \\ \|\hat{\mathcal{W}}_\sigma - \mathcal{W}_\sigma\|_F &= O_P(n^{-\frac{1}{2}}). \end{aligned} \quad (24)$$

By (24), we have that (21) holds with $\delta = O_P(n^{-\frac{1}{2}})$. Hence, by Lemma 3, the eigenpairs of $\hat{\mathcal{W}}_\sigma$ provide a candidate set \hat{V}_σ that contains d vectors that are $O_P(n^{-\frac{1}{2}})$ close to the columns of W^\dagger . In addition, any irrelevant candidate is $\Omega_P(1)$ far away from W^\dagger . As we show next, these properties ensure that with high probability the d relevant candidates can be identified in \hat{V}_σ .

Filtering Given the candidate set \hat{V}_σ computed in the first step, our goal now is to find a set $\tilde{V}_\sigma \subseteq \hat{V}_\sigma$ of d vectors that accurately estimate the d columns of W^\dagger . To simplify the theoretical analysis, we assume the filtering step is done using a sample X of size n that is independent of \tilde{V}_σ . This can be achieved by first splitting a given sample of size $2n$ into two sets of size n , one for each step.

Recall that for \boldsymbol{x} from model (1) and any $\boldsymbol{v} \in \mathbb{R}^m$,

$$\boldsymbol{v}^\top \boldsymbol{x} = \boldsymbol{v}^\top W^\top \boldsymbol{h} + \sigma \boldsymbol{v}^\top \boldsymbol{\xi}. \quad (25)$$

Obviously, when $\sigma > 0$, the filtering procedure in (16) for the noiseless case is inadequate, as typically no $\boldsymbol{v}^* \in W^\dagger$ will exactly satisfy $\boldsymbol{v}^{*\top} X \in \{0, 1\}^n$. Nevertheless, we expect that for a sufficiently small noise level σ , any $\boldsymbol{v} \in \hat{V}_\sigma$ that is close to some $\boldsymbol{v}^* \in W^\dagger$ will result in $\boldsymbol{v}^\top X$ that is close to being binary, while any \boldsymbol{v} sufficiently far from W^\dagger will result in $\boldsymbol{v}^\top X$ that is far from being binary. A natural measure for how $\boldsymbol{v}^\top X$ is ‘‘far from being binary’’, similar to the one used for filtering in Slawski et al. (2013), is simply its deviation from its binary rounding,

$$\min_{\boldsymbol{b} \in \{0, 1\}^n} \frac{\|\boldsymbol{v}^\top X - \boldsymbol{b}\|^2}{n \|\boldsymbol{v}\|^2}. \quad (26)$$

Eq. (26) works extremely well for small σ , but fails for high noise levels. Here we instead propose a filtering procedure based on the classical Kolmogorov-Smirnov goodness of fit

test (Lehmann & Romano, 2006). As we show below, this approach gives consistent estimates of W for any $\sigma > 0$.

Before describing the test, we first introduce the probabilistic analogue of the rigidity condition (6). For any $\mathbf{u} \in \mathbb{R}^d$, define its corresponding expected binary rounding error,

$$r(\mathbf{u}) = \mathbb{E}_{\mathbf{h} \sim P_{\mathbf{h}}} \left[\min_{b \in \{0,1\}} (\mathbf{u}^\top \mathbf{h} - b)^2 \right].$$

Clearly, $r(\mathbf{0}) = 0$ and $r(\mathbf{e}_i) = 0$ for all $i \in [d]$. We pose the following *expected rigidity* condition: for all $\mathbf{u} \neq \mathbf{0}$,

$$r(\mathbf{u}) = 0 \quad \Leftrightarrow \quad \mathbf{u} \in \{\mathbf{e}_i\}_{i=1}^d. \quad (27)$$

Analogously to the deterministic rigidity condition in (6), condition (27) is satisfied, for example, when $P_{\mathbf{h}}(\mathbf{e}_i) > 0$ and $P_{\mathbf{h}}(\mathbf{e}_i + \mathbf{e}_j) > 0$ for all $i \neq j \in [d]$.

To introduce our filtering test, recall that under model (1), $\boldsymbol{\xi} \sim \mathcal{N}(\mathbf{0}, I_m)$. Hence, for any fixed \mathbf{v} , the random variable $\mathbf{v}^\top \mathbf{x}$ in (25) is distributed according to the following univariate Gaussian mixture model (GMM),

$$\mathbf{v}^\top \mathbf{x} \sim \sum_{\mathbf{h} \in \{0,1\}^d} P_{\mathbf{h}}(\mathbf{h}) \cdot \mathcal{N}(\mathbf{v}^\top W^\top \mathbf{h}, \sigma^2 \|\mathbf{v}\|^2). \quad (28)$$

Denote the cumulative distribution function of $\mathbf{v}^\top \mathbf{x}$ by $F_{\mathbf{v}}$. For general \mathbf{v} , this mixture may have up to 2^d distinct components. However, for $\mathbf{v}^* \in W^\dagger$, it reduces to a mixture of *two* components with means at 0 and 1. More precisely, for any candidate \mathbf{v} with corresponding eigenvalue $\lambda(\mathbf{v}) \geq 1$, define the GMM with two components

$$(1 - \frac{1}{\lambda(\mathbf{v})^2}) \cdot \mathcal{N}(0, \sigma^2 \|\mathbf{v}\|^2) + \frac{1}{\lambda(\mathbf{v})^2} \cdot \mathcal{N}(1, \sigma^2 \|\mathbf{v}\|^2). \quad (29)$$

Denote its cumulative distribution function by $G_{\mathbf{v}}$. The following lemma shows that under condition (27), $G_{\mathbf{v}}$ fully characterizes the columns of W^\dagger .

Lemma 4. *Let $K_\sigma, \mathcal{W}_\sigma$ be the population quantities in (19) and let V_σ be the set of population candidates as computed from the eigenpairs of \mathcal{W}_σ . If conditions (I)-(II) and the expected rigidity condition (27) hold, then for any $\mathbf{v} \in V_\sigma$ with corresponding eigenvalue $\lambda(\mathbf{v})$,*

$$F_{\mathbf{v}} = G_{\mathbf{v}} \quad \Leftrightarrow \quad \mathbf{v} \in W^\dagger.$$

Given the empirical candidate set \hat{V}_σ , Lemma 4 suggests ranking all $\hat{\mathbf{v}} \in \hat{V}_\sigma$ according to their goodness of fit to $G_{\hat{\mathbf{v}}}$ and taking the d candidates with the best fit. More precisely, given a sample $X = [\mathbf{x}_1, \dots, \mathbf{x}_n]$ that is independent of \hat{V}_σ , for each candidate $\hat{\mathbf{v}} \in \hat{V}_\sigma$ we compute the empirical cumulative distribution function, $\hat{F}_{\hat{\mathbf{v}}}(t) = \frac{1}{n} \sum_{j=1}^n \mathbb{1}\{\hat{\mathbf{v}}^\top \mathbf{x}_j \leq t\}$, $t \in \mathbb{R}$, and calculate its Kolmogorov-Smirnov score

$$\Delta_n(\hat{\mathbf{v}}) = \sup_{t \in \mathbb{R}} |\hat{F}_{\hat{\mathbf{v}}}(t) - G_{\hat{\mathbf{v}}}(t)|. \quad (30)$$

Algorithm 1 Estimate W when $\sigma > 0$ and $n < \infty$

Input: sample matrix $X \in \mathbb{R}^{m \times n}$ and $0 < \tau \ll 1$

- 1: estimate number of hidden units d and noise level σ^2
- 2: compute empirical moments $\hat{\boldsymbol{\mu}}, \hat{M}$ and $\hat{\mathcal{M}}$ and plugin moments \hat{M}_σ and $\hat{\mathcal{M}}_\sigma$ of (17)
- 3: compute \hat{K}_σ such that $\hat{K}_\sigma^\top \hat{M}_\sigma \hat{K}_\sigma = I_d$
- 4: construct $\hat{\mathcal{W}}_\sigma = \hat{\mathcal{M}}_\sigma(\hat{K}_\sigma, \hat{K}_\sigma, \hat{K}_\sigma)$
- 5: compute the set \hat{U}_σ of eigenpairs of $\hat{\mathcal{W}}_\sigma$
- 6: compute the candidate set \hat{V}_σ in (20)
- 7: for each $\hat{\mathbf{v}} \in \hat{V}_\sigma$ compute its KS score $\Delta_n(\hat{\mathbf{v}})$ in (30)
- 8: select $\bar{V}_\sigma \subseteq \hat{V}_\sigma$ of d vectors with smallest $\Delta_n(\hat{\mathbf{v}})$
- 9: **return** the pseudo-inverse $\hat{W} = \bar{V}_\sigma^\dagger$

Our estimator $\bar{V}_\sigma \subseteq \hat{V}_\sigma$ for W^\dagger is then the set of d vectors with the smallest scores $\Delta_n(\hat{\mathbf{v}})$. The estimator for W is the pseudo-inverse, $\hat{W} = \bar{V}_\sigma^\dagger$.

The following lemma shows that for sufficiently large n , $\Delta_n(\hat{\mathbf{v}})$ accurately distinguishes between $\hat{\mathbf{v}} \in \hat{V}_\sigma$ that are close to the columns of W^\dagger from these that are not.

Lemma 5. *Let $\mathbf{v}^* \in W^\dagger$ and $\hat{\mathbf{v}}_{(1)}, \hat{\mathbf{v}}_{(2)}, \dots$ a sequence of random vectors such that $\|\hat{\mathbf{v}}_{(n)} - \mathbf{v}^*\| = O_P(n^{-\frac{1}{2}})$. Then, $\Delta_n(\hat{\mathbf{v}}_{(n)}) = o_P(1)$. In contrast, if $\min_{\mathbf{v}^* \in W^\dagger} \|\hat{\mathbf{v}}_{(n)} - \mathbf{v}^*\| = \Omega_P(1)$, then $\Delta_n(\hat{\mathbf{v}}_{(n)}) = \Omega_P(1)$, provided the expected rigidity condition (27) holds.*

Lemma 5 follows from classical and well studied properties of the Kolmogorov-Smirnov test, see for example Lehmann & Romano (2006); Billingsley (2013).

Algorithm 1 summarizes our method for estimating W in the general case where $\sigma > 0$ and $n < \infty$. The following theorem establishes its consistency.

Theorem 2. *Let $\mathbf{x}_1, \dots, \mathbf{x}_n$ be n i.i.d. samples from model (1). If conditions (I)-(II) and the expected rigidity condition (27) hold, then the estimator \hat{W} computed by Algorithm 1 is consistent, achieving the parametric rate,*

$$\hat{W} = W + O_P(n^{-\frac{1}{2}}).$$

Runtime The runtime of Algorithm 1 is composed of three main parts. First, $O(nm^3)$ operations are needed to compute all the relevant moments from the data and to construct the $d \times d \times d$ whitened tensor $\hat{\mathcal{W}}_\sigma$. The most time consuming task is computing the eigenpairs of $\hat{\mathcal{W}}_\sigma$, which can be done by either the homotopy method or O-NCM. Currently, no runtime guarantees are available for either of these methods. In practice, since there are $O(2^d)$ eigenpairs, these methods spend $O(2^d \cdot \text{poly}(d))$ operations in total. Finally, since there are $O(2^d)$ candidates and each KS test takes $O(dn)$ operations (Gonzalez et al., 1977), the filtering procedure runtime is $O(d2^d n)$.

Power-stability and orthogonal decomposition The exponential runtime of our algorithm stems from the fact that the set U_N of Newton-stable eigenpairs of \mathcal{W}_σ is typically $O(2^d)$. However, in some cases, the set U^* of d relevant eigenpairs has additional structure so that a smaller candidate set may be computed instead of U_N . Consider the subset $U_P \subseteq U_N$ of *power-stable* eigenpairs of \mathcal{W}_σ :

Definition 3. An eigenpair (\mathbf{u}, λ) is *power-stable* if its projected Jacobian $J_p(\mathbf{u})$ is either positive or negative definite.

Typically, the number of power-stable eigenpairs is significantly smaller than the number of Newton-stable eigenpairs.² In addition, U_P can be computed by the shifted higher-order power method (Kolda & Mayo, 2011; 2014).

Similarly to Lemma 2, one can show that U_P is guaranteed to contain U^* whenever the following stronger version of condition (II) holds: for all $(\mathbf{u}_i, \lambda_i) \in U^*$, the matrix

$$(WKL\mathbf{u}_i)^\top (2\mathcal{C}(I, I, \mathbf{e}_i) - \mathcal{C})(WKL\mathbf{u}_i) \quad (31)$$

is either positive-definite or negative-definite.

As an example, consider the case where P_h has the support $\mathbf{h} \in I_d$. Then model (1) corresponds to a GMM with d spherical components with linearly independent means. In this case, both \mathcal{C} and \mathcal{C} are diagonal with \mathbf{p} on their diagonal. Thus, the matrices in (31) take the form $-L_{\mathbf{e}_i}^\top \text{diag}(\mathbf{p})L_{\mathbf{e}_i}$, which are all negative-definite when $\mathbf{p} > 0$. In fact, in this case, \mathcal{W}_σ has an orthogonal CP decomposition and the d orthogonal eigenpairs in U^* are the *only* negative-definite power-stable eigenpairs of \mathcal{W}_σ (Anandkumar et al., 2014). Similarly, when P_h is a product distribution, the same orthogonal structure appears if the *centered* moments of \mathbf{x} are used instead of M and \mathcal{M} . As shown in Anandkumar et al. (2014), the power method, accompanied with a deflation procedure, decomposes an orthogonal tensor in polynomial time, thus implying an efficient algorithm in these cases. However, under the much weaker conditions we pose on P_h , the relevant eigenpairs in U^* are not necessarily power-stable and the CP decomposition of \mathcal{W}_σ does not necessarily include U^* .

5. Experiments

We demonstrate our method in three scenarios: (I) simulations from the exact binary model (1); (II) learning a common population genetic admixture model; (III) learning the proportion matrix of a cell mixture from DNA methylation levels. Due to lack of space, (III) is deferred to Appendix N. Code to reproduce the simulation results can be found at <https://github.com/arJaffe/BinaryLatentVariables>.

²We currently do not know whether the number of power-stable eigenpairs of a generic tensor is polynomial or exponential in d .

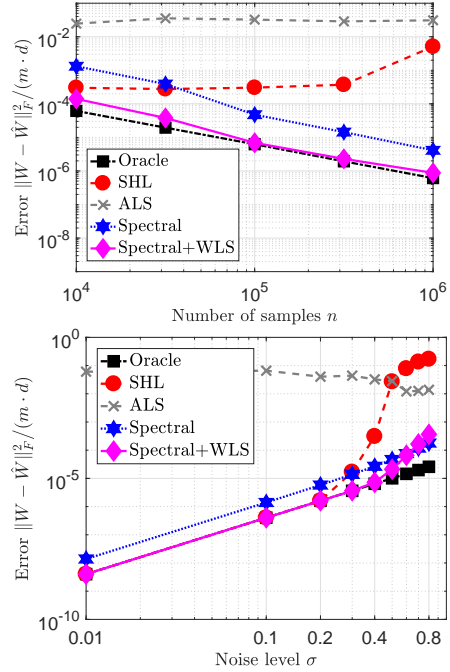


Figure 2. Upper panel: Error in W vs. sample size n with $\sigma = 0.4$. lower panel: Error in W vs. noise level σ with $n = 10^5$.

5.1. Simulations

We generated n samples from model (1) with $d = 6$ hidden units, $m = 30$ observable features, and Gaussian noise $\xi \sim \mathcal{N}(\mathbf{0}, I_m)$. The m columns of W were drawn uniformly from the unit sphere \mathbb{S}_{d-1} . Fixing a mean vector $\mathbf{a} \in \mathbb{R}^d$ and a covariance matrix $R \in \mathbb{R}^{d \times d}$, each hidden vector \mathbf{h} was generated independently by first drawing $\mathbf{r} \sim \mathcal{N}(\mathbf{a}, R)$ and then taking its binary rounding.

Figure 2 shows the error, in Frobenius norm, averaged over 50 independent realizations of X as a function of n (upper panel) and σ (lower panel) for 5 methods: (i) our spectral approach, Algorithm 1 (Spectral); (ii) Algorithm 1 followed by a single weighted least squares step (Appendix K) (Spectral+WLS); (iii) SHL, the matrix decomposition method of Slawski et al. (2013)³; (iv) ALS with a random initialization (Appendix L); and (v) an oracle estimator that is given the exact matrix H and computes W via least squares.

As one can see, as opposed to SHL, our method is consistent for $\sigma > 0$ and achieves an error rate $O(n^{-\frac{1}{2}})$ corresponding to a slope of -1 in the upper panel of Fig. 2. In addition, as seen in the lower panel of Fig. 2, at low levels of noise our method is comparable to SHL, whereas at high levels it is far more accurate. Finally, adding a weighted least squares step reduces the error for low noise levels, but increases the error

³Code from <https://sites.google.com/site/slawskimartin/code>. For each realization X , we made 50 runs of SHL and chose H, W minimizing $\|X - W^\top H\|_F$.

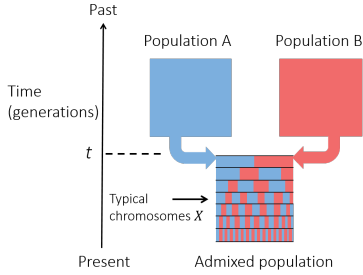


Figure 3. Illustration of the admixture model.

for high noise levels. A comparison between the runtime of SHL and the spectral method appears in Appendix I.

5.2. Population genetic admixture

We present an application of our method to a fundamental problem in population genetics, known as admixture (see Fig. 3). Admixture refers to the mixing of $d \geq 2$ ancestral populations that were long separated, e.g., due to geographical or cultural barriers (Pritchard et al., 2000; Alexander et al., 2009; Li et al., 2008). The observed data X is an $m \times n$ matrix where m is the number of modern “admixed” individuals and n is the number of relevant locations in their DNA, known as SNPs. Each SNP corresponds to two alleles and individuals may have different alleles. Fixing a reference allele for each location, X_{ij} takes values in $\{0, \frac{1}{2}, 1\}$ according to the number of reference alleles appearing in the genotype of individual $i \in [m]$ at locus $j \in [n]$.

Given the genotypes X , an important problem in population genetics is to estimate the following two quantities. The allele frequency matrix $H \in [0, 1]^{d \times n}$ whose entry H_{kj} is the frequency of the reference allele at locus $j \in [n]$ in ancestral population $k \in [d]$; and the admixture proportion matrix $W \in [0, 1]^{d \times m}$ whose columns sum to 1 and its entry W_{ki} is the proportion of individual i ’s genome that was inherited from population k .

A common model for X in terms of W and H is to assume that the number of alleles $2X_{ij} \in \{0, 1, 2\}$ is the sum of two i.i.d. Bernoulli random variables with success probability $F_{ij} = \sum_{k=1}^d W_{ki} H_{kj}$, namely, $X_{ij} | H \sim \frac{1}{2} \cdot \text{Binomial}(2, F_{ij})$. Note that under this model

$$\mathbb{E}[X|H] = F = W^\top H. \quad (32)$$

Although (32) has similar form to model (1), there are two main differences; the noise is not normally distributed and the matrix H is non-binary. Yet, the binary model (1) serves as a good approximation whenever various alleles are rare in some populations but abundant in others. Specifically, for ancestral populations that have been long separated, some alleles may become *fixed* in one population (i.e., reach frequency of 1) while being totally absent in others.

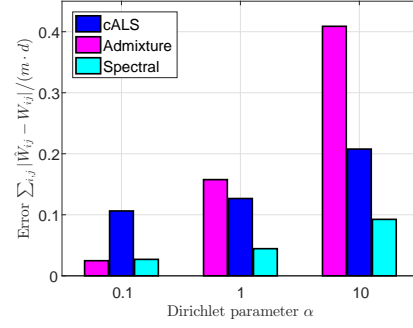


Figure 4. Average absolute error in \hat{W} vs. Dirichlet parameter α for $d = 3$ ancestral populations and $m = 50$ admixed individuals.

Simulating genetic admixture We followed a standard simulation scheme applied, for example, in Xue et al. (2017); Gravel (2012); Price et al. (2009). First, using SCRM (Staab et al., 2015), we simulated $d = 3$ ancestral populations separated for 4000 generations and generated the genomes of 40 individuals for each. H was then computed as the frequency of the reference alleles in each population. Next, the columns of W were sampled from a symmetric Dirichlet distribution with parameter $\alpha \geq 0$. Finally, the genomes of $m = 50$ admixed individuals were generated as mosaics of genomic segments of individuals from the ancestral populations with proportions W . The mosaic nature of the admixed genomes is an important realistic detail, due to the *linkage* (correlation) between SNPs (Xue et al., 2017). A detailed description is in Appendix M.

We compare our algorithm to two methods. The first is Admixture (Alexander et al., 2009), one of the most widely used algorithms in population genetics, which aims to maximize the likelihood of X . The second is the recently proposed spectral method ALStructure (Cabreros & Storey, 2017), where an estimation of $\text{span}(W^\top)$ via Chen & Storey (2015) is followed by constrained ALS iterations of W and H . For our method, two modifications are needed for Algorithm 1. First, since the distribution of $X_{ij} - w_i^T h_j$ is not Gaussian, the corrected moments $\hat{M}_\sigma, \hat{\mathcal{M}}_\sigma$ as calculated by (17) do not satisfy (18). Instead, we implemented a matrix completion algorithm derived in (Jain & Oh, 2014) for a similar setup, see Appendix J for more details. In addition, the filtering process described in Section 4 is no longer valid. However, as d is relatively small, we performed exhaustive search over all candidate subsets of size d and choose the one that maximized the likelihood.

Figure 4 compares the results of the 3 methods for $\alpha = 0.1, 1, 10$. The spectral method outperforms Admixture and ALStructure for $\alpha = 1, 10$ and performs similarly to Admixture for $\alpha = 0.1$.

Acknowledgements This research was funded in part by NIH Grant 1R01HG008383-01A1.

References

- Alexander, D., Novembre, J., and Lange, K. Fast model-based estimation of ancestry in unrelated individuals. *Genome research*, 19(9):1655–1664, 2009.
- Anandkumar, A., Foster, D., Hsu, D., Kakade, S., and Liu, Y. A spectral algorithm for latent dirichlet allocation. In *Advances in Neural Information Processing Systems*, pp. 917–925, 2012a.
- Anandkumar, A., Hsu, D., and Kakade, S. A method of moments for mixture models and hidden markov models. In *COLT*, volume 1, pp. 4, 2012b.
- Anandkumar, A., Ge, R., Hsu, D., Kakade, S., and Telgarsky, M. Tensor decompositions for learning latent variable models. *Journal of Machine Learning Research*, 15(1): 2773–2832, 2014.
- Arabshahi, F. and Anandkumar, A. Spectral methods for correlated topic models. In *Artificial Intelligence and Statistics*, pp. 1439–1447, 2017.
- Arora, S., Ge, R., Kannan, R., and Moitra, A. Computing a nonnegative matrix factorization—provably. In *Proceedings of the forty-fourth annual ACM symposium on Theory of computing*, pp. 145–162. ACM, 2012.
- Baadel, S., Thabtah, F., and Lu, J. Overlapping clustering: A review. In *SAI Computing Conference (SAI), 2016*, pp. 233–237. IEEE, 2016.
- Banerjee, A., Krumpelman, C., Ghosh, J., Basu, S., and Mooney, R. Model-based overlapping clustering. In *Proceedings of the eleventh ACM SIGKDD international conference on Knowledge discovery in data mining*, pp. 532–537. ACM, 2005.
- Becker, E., Robisson, B., Chapple, C., Guénoche, A., and Brun, C. Multifunctional proteins revealed by overlapping clustering in protein interaction network. *Bioinformatics*, 28(1):84–90, 2011.
- Billingsley, P. *Convergence of probability measures*. John Wiley & Sons, 2013.
- Cabreros, I. and Storey, J. A nonparametric estimator of population structure unifying admixture models and principal components analysis. *bioRxiv*, pp. 240812, 2017.
- Cartwright, D. and Sturmfels, B. The number of eigenvalues of a tensor. *Linear algebra and its applications*, 438(2): 942–952, 2013.
- Chen, L., Han, L., and Zhou, L. Computing tensor eigenvalues via homotopy methods. *SIAM Journal on Matrix Analysis and Applications*, 37(1):290–319, 2016.
- Chen, X. and Storey, J. Consistent estimation of low-dimensional latent structure in high-dimensional data. *arXiv preprint arXiv:1510.03497*, 2015.
- Cho, Ilin, and Raiko. Improved learning of Gaussian-Bernoulli restricted Boltzmann machines. *Artificial Neural Networks and Machine Learning—ICANN 2011*, pp. 10–17, 2011.
- Cui, C., Dai, Y., and Nie, J. All real eigenvalues of symmetric tensors. *SIAM Journal on Matrix Analysis and Applications*, 35(4):1582–1601, 2014.
- Donoho, D. and Stodden, V. When does non-negative matrix factorization give a correct decomposition into parts? In *Advances in neural information processing systems*, pp. 1141–1148, 2004.
- Frieze, A., Jerrum, M., and Kannan, R. Learning linear transformations. In *Foundations of Computer Science, 1996. Proceedings., 37th Annual Symposium on*, pp. 359–368. IEEE, 1996.
- Golub, G. and Van Loan, C. *Matrix computations*, volume 3. JHU Press, 2012.
- Gonzalez, T., Sahni, S., and Franta, W. An efficient algorithm for the Kolmogorov-Smirnov and Lilliefors tests. *ACM Transactions on Mathematical Software (TOMS)*, 3(1):60–64, 1977.
- Gravel, S. Population genetics models of local ancestry. *Genetics*, 191(2):607–619, 2012.
- Guo, C., Lin, W., and Liu, C. A modified Newton iteration for finding nonnegative Z-eigenpairs of a nonnegative tensor. *arXiv preprint arXiv:1705.07487*, 2017.
- Hillar, C. and Lim, L. Most tensor problems are np-hard. *Journal of the ACM (JACM)*, 60(6):45, 2013.
- Hinton, G. A practical guide to training restricted Boltzmann machines. *Momentum*, 9(1):926, 2010.
- Hinton, G. and Salakhutdinov, R. Reducing the dimensionality of data with neural networks. *science*, 313(5786): 504–507, 2006.
- Houseman, E., Accomando, P., Koestler, D., Christensen, B., Marsit, C., Nelson, H., Wiencke, J., and Kelsey, K. Dna methylation arrays as surrogate measures of cell mixture distribution. *BMC bioinformatics*, 13(1):86, 2012.
- Hsu, D. and Kakade, S. Learning mixtures of spherical gaussians: moment methods and spectral decompositions. In *Proceedings of the 4th conference on Innovations in Theoretical Computer Science*, pp. 11–20. ACM, 2013.

- Hubbard, J. and Hubbard, B. *Vector calculus, linear algebra, and differential forms: a unified approach*. Matrix Editions, 2015.
- Hyvärinen, A., Karhunen, J., and Oja, E. *Independent component analysis*, volume 46. John Wiley & Sons, 2004.
- Jaffe, A., Weiss, R., and Nadler, B. Newton correction methods for computing real eigenpairs of symmetric tensors. *arXiv preprint arXiv:1706.02132*, 2017.
- Jain, P. and Oh, S. Learning mixtures of discrete product distributions using spectral decompositions. In *Conference on Learning Theory*, pp. 824–856, 2014.
- Kolda, T. and Mayo, J. Shifted power method for computing tensor eigenpairs. *SIAM Journal on Matrix Analysis and Applications*, 32(4):1095–1124, 2011.
- Kolda, T. and Mayo, J. An adaptive shifted power method for computing generalized tensor eigenpairs. *SIAM Journal on Matrix Analysis and Applications*, 35(4):1563–1581, 2014.
- Kritchman, S. and Nadler, B. Non-parametric detection of the number of signals: Hypothesis testing and random matrix theory. *IEEE Transactions on Signal Processing*, 57(10):3930–3941, 2009.
- Lehmann, E. and Romano, J. *Testing statistical hypotheses*. Springer Science & Business Media, 2006.
- Li, J., Absher, D., Tang, H., Southwick, A., Casto, A., Ramachandran, S., Cann, H., Barsh, G., Feldman, M., Cavalli-Sforza, L., et al. Worldwide human relationships inferred from genome-wide patterns of variation. *science*, 319(5866):1100–1104, 2008.
- Lim, L. Singular values and eigenvalues of tensors: a variational approach. In *Computational Advances in Multi-Sensor Adaptive Processing, 2005 1st IEEE International Workshop on*, pp. 129–132. IEEE, 2005.
- Melchior, J., Wang, N., and Wiskott, L. Gaussian-binary restricted Boltzmann machines for modeling natural image statistics. *PLoS one*, 12(2):e0171015, 2017.
- Miettinen, P. and Vreeken, J. Mdl4bmf: Minimum description length for boolean matrix factorization. *ACM Transactions on Knowledge Discovery from Data (TKDD)*, 8(4):18, 2014.
- Moitra, A. and Valiant, G. Settling the polynomial learnability of mixtures of gaussians. In *Foundations of Computer Science (FOCS), 2010 51st Annual IEEE Symposium on*, pp. 93–102. IEEE, 2010.
- Mossel, E. and Roch, S. Learning nonsingular phylogenies and hidden markov models. In *Proceedings of the thirty-seventh annual ACM symposium on Theory of computing*, pp. 366–375. ACM, 2005.
- Price, A., Tandon, A., Patterson, N., Barnes, K., Rafaels, N., Ruczinski, I., Beaty, T., Mathias, R., Reich, D., and Myers, S. Sensitive detection of chromosomal segments of distinct ancestry in admixed populations. *PLoS genetics*, 5(6):e1000519, 2009.
- Pritchard, J., Stephens, M., and Donnelly, P. Inference of population structure using multilocus genotype data. *Genetics*, 155(2):945–959, 2000.
- Qi, L. Eigenvalues of a real supersymmetric tensor. *Journal of Symbolic Computation*, 40(6):1302–1324, 2005.
- Regalia, P. and Kofidis, E. Monotonic convergence of fixed-point algorithms for ica. *IEEE Transactions on Neural Networks*, 14(4):943–949, 2003.
- Segal, E., Battle, A., and Koller, D. Decomposing gene expression into cellular processes. In *Proceedings of the Pacific Symposium on Biocomputing*, volume 8, pp. 89–100, 2002.
- Shalvi, O. and Weinstein, E. Super-exponential methods for blind deconvolution. *IEEE Transactions on Information Theory*, 39(2):504–519, 1993.
- Slawski, M., Hein, M., and Lutsik, P. Matrix factorization with binary components. In *Advances in Neural Information Processing Systems*, pp. 3210–3218, 2013.
- Staab, Zhu, Metzler, and Lunter. scrm: efficiently simulating long sequences using the approximated coalescent with recombination. *Bioinformatics*, 31(10):1680–1682, 2015.
- Taylor, G., Hinton, G., and Roweis, S. Modeling human motion using binary latent variables. In *Advances in neural information processing systems*, pp. 1345–1352, 2007.
- Van der Veen, A.-J. Analytical method for blind binary signal separation. *IEEE Transactions on Signal Processing*, 45(4):1078–1082, 1997.
- Wang, N., Melchior, J., and Wiskott, L. An analysis of gaussian-binary restricted Boltzmann machines for natural images. In *ESANN*, 2012.
- Xue, J., Lencz, T., Darvasi, A., Peer, I., and Carmi, S. The time and place of european admixture in ashkenazi jewish history. *PLoS genetics*, 13(4):e1006644, 2017.

A. Proof of Proposition 1

Uniqueness of the factorization readily follows from (7) so we proceed to prove (7). First note that $\text{span}(X) = \text{span}(W^\top) = \text{span}(W^\dagger)$. Since W is full rank, we have $WW^\dagger = I_d$. Hence,

$$(W^\dagger)^\top X = (WW^\dagger)^\top H = H \in \{0, 1\}^{d \times n}.$$

So any $\mathbf{v}^* \in W^\dagger$ satisfies the binary constraint $\mathbf{v}^{*\top} X \in \{0, 1\}^n$. For the other direction, let $\mathbf{v} \in \text{span}(X) \setminus \{0\}$ be such that $\mathbf{v}^\top X \in \{0, 1\}^n$. Since $\mathbf{v}^\top X = (W\mathbf{v})^\top H$, the rigidity condition (6) implies $W\mathbf{v} \in \{\mathbf{e}_i\}_{i=1}^d$. Since W is full rank and $\mathbf{v} \in \text{span}(W^\dagger)$, \mathbf{v} must be a column of W^\dagger .

B. Proof of Lemma 1

Since the vector \mathbf{h} is binary, its second and third order moments are related as follows. For all $i, j \in [d]$,

$$C_{iij} = C_{iji} = C_{jii} = \mathbb{E}[h_i^2 h_j] = \mathbb{E}[h_i h_j] = C_{ij}. \quad (33)$$

Since W is full rank, $WW^\dagger = I_d$. Hence, applying W^\dagger multi-linearly on the moment equations in (10) we obtain

$$\begin{aligned} C &= (W^\dagger)^\top M W^\dagger, \\ \mathcal{C} &= \mathcal{M}(W^\dagger, W^\dagger, W^\dagger). \end{aligned}$$

Thus, the equality in (33) is equivalent to

$$[\mathcal{M}(W^\dagger, W^\dagger, W^\dagger)]_{iij} = [(W^\dagger)^\top M W^\dagger]_{ij}. \quad (34)$$

Let $Y^* \in \mathbb{R}^{d \times d}$ be the full rank matrix that satisfies $W^\dagger = KY^*$ where K is the whitening matrix in (11). Then,

$$\begin{aligned} \mathcal{M}(W^\dagger, W^\dagger, W^\dagger) &= \mathcal{M}(KY^*, KY^*, KY^*) \\ &= \mathcal{W}(Y^*, Y^*, Y^*) \end{aligned}$$

where \mathcal{W} is the whitened tensor in (12). Similarly, by (11),

$$(W^\dagger)^\top M W^\dagger = (Y^*)^\top (K^\top M K) (Y^*) = (Y^*)^\top Y^*.$$

Inserting these into (34), the matrix Y^* must satisfy

$$[\mathcal{W}(Y^*, Y^*, Y^*)]_{iij} = [(Y^*)^\top (Y^*)]_{ij}, \quad \forall i, j \in [d]. \quad (35)$$

The following lemma, proved in Appendix H, shows that Eq. (35) is nothing but a tensor eigen-problem. Specifically, the columns of Y^* , up to scaling, are eigenvectors of \mathcal{W} .

Lemma 6. *Let $\mathcal{W} \in \mathbb{R}^{d \times d \times d}$ be an arbitrary symmetric tensor. Then, a matrix $Y = [\mathbf{y}_1, \dots, \mathbf{y}_d] \in \mathbb{R}^{d \times d}$ of rank d satisfies (35) if and only if for all $k \in [d]$, $\mathbf{y}_k = \mathbf{u}_k / \lambda_k$, where $(\mathbf{u}_k, \lambda_k)_{k=1}^d$ are d eigenpairs of \mathcal{W} with linearly independent $\{\mathbf{u}_k\}_{k=1}^d$.*

By Lemma 6, the set of scaled eigenpairs $\{\mathbf{y} = \mathbf{u}/\lambda\}$ of \mathcal{W} is guaranteed to contain the d columns of Y^* . Since $W^\dagger = KY^*$, the set $\{K\mathbf{y}\}$ is guaranteed to contain W^\dagger .

To show that each $\mathbf{y} = \mathbf{u}/\lambda \in Y^*$ has $\lambda \geq 1$, note that the vector $K\mathbf{y}$ is a column of W^\dagger , so $WK\mathbf{y} = \mathbf{e}_i$ for some $i \in [d]$. Hence, by the definition of the whitened tensor (12) and the moment equation (10),

$$\begin{aligned} \mathcal{W}(\mathbf{y}, \mathbf{y}, \mathbf{y}) &= \mathcal{M}(K\mathbf{y}, K\mathbf{y}, K\mathbf{y}) \\ &= \mathcal{C}(WK\mathbf{y}, WK\mathbf{y}, WK\mathbf{y}) \\ &= \mathcal{C}(\mathbf{e}_i, \mathbf{e}_i, \mathbf{e}_i) = C_{iii} = \mathbb{E}[h_i] \leq 1. \end{aligned}$$

On the other hand, since (\mathbf{u}, λ) is an eigenpair of \mathcal{W} with eigenvalue $\lambda = \mathcal{W}(\mathbf{u}, \mathbf{u}, \mathbf{u})$,

$$\mathcal{W}(\mathbf{y}, \mathbf{y}, \mathbf{y}) = \frac{1}{\lambda^3} \mathcal{W}(\mathbf{u}, \mathbf{u}, \mathbf{u}) = \frac{1}{\lambda^2}.$$

By convention, $\lambda \geq 0$. Hence,

$$\lambda = 1/\sqrt{\mathbb{E}[h_i]} \geq 1,$$

concluding the proof.

C. Recovery algorithm - noiseless case

Algorithm 2 Recover W when $\sigma = 0$

Input: sample matrix X

- 1: estimate second and third order moments M, \mathcal{M}
 - 2: set $d = \text{rank}(M)$
 - 3: compute $K \subseteq \text{span}(M)$ such that $K^\top M K = I_d$
 - 4: compute whitened tensor $\mathcal{W} = \mathcal{M}(K, K, K)$
 - 5: compute the set U of eigenpairs of \mathcal{W}
 - 6: compute the candidate set V in (14)
 - 7: filter $\bar{V} = \{\mathbf{v} \in V : \mathbf{v}^\top X \in \{0, 1\}^n\}$
 - 8: **return** the pseudo-inverse $W = \bar{V}^\dagger$
-

D. Proof of Lemma 2

Let $(\mathbf{u}, \lambda) \in U^*$ be an eigenpair of \mathcal{W} such that $\mathbf{v}^* = K\mathbf{u}/\lambda \in W^\dagger$. To show Newton-stability we need to show that under conditions (I)-(II) the projected Jacobian matrix $J_p(\mathbf{u}) = L_u^\top \nabla \mathbf{g}(\mathbf{u}) L_u$ in (4) is full rank $d - 1$.

The Jacobian matrix $\nabla \mathbf{g}(\mathbf{u})$ is

$$\begin{aligned} \nabla \mathbf{g}(\mathbf{u}) &= 2\mathcal{W}(I, I, \mathbf{u}) - 3\mathbf{u}\mathcal{W}(I, \mathbf{u}, \mathbf{u})^\top \\ &\quad - \mathcal{W}(\mathbf{u}, \mathbf{u}, \mathbf{u}) I_d \\ &= 2\mathcal{W}(I, I, \mathbf{u}) - 3\lambda \mathbf{u}\mathbf{u}^\top - \lambda I_d. \end{aligned} \quad (36)$$

Since $L_u^\top \mathbf{u} = \mathbf{0}$, the second term in (36) does not contribute to $J_p(\mathbf{u})$. For the first term in (36), by (12) and (10),

$$\mathcal{W}(I, I, \mathbf{u}) = \mathcal{M}(K, K, K\mathbf{u}) = \mathcal{C}(WK, WK, WK\mathbf{u}).$$

Since $\mathbf{v}^* = K\mathbf{u}/\lambda$ is a column of W^\dagger , $WK\mathbf{u} = \lambda \mathbf{e}_i$ for some $i \in [d]$. Thus,

$$\begin{aligned} \mathcal{W}(I, I, \mathbf{u}) &= \lambda \mathcal{C}(WK, WK, \mathbf{e}_i) \\ &= \lambda K^\top W^\top \mathcal{C}(I, I, \mathbf{e}_i) WK. \end{aligned}$$

For the third term in (36), by the definition of K in (11),

$$I_d = K^\top MK = K^\top W^\top CWK.$$

Putting the last two equalities in (36) and applying the projection L_u we obtain

$$\begin{aligned} J_p(\mathbf{u}) &= L_u^\top \nabla g(\mathbf{u}) L_u \\ &= \lambda L_u^\top K^\top W^\top (2C(I, I, e_i) - C) W K L_u. \end{aligned}$$

Since $\lambda \geq 1$ and W and K are full rank, condition (II) implies that $J_p(\mathbf{u})$ is full rank as well. Thus, (\mathbf{u}, λ) is a Newton-stable eigenpair of \mathcal{W} .

E. Proof of Lemma 3

Lemma 3 follows from the following lemma which establishes the stability of Newton-stable eigenpairs of a tensor \mathcal{W} to small perturbations $\tilde{\mathcal{W}} = \mathcal{W} + \Delta\mathcal{W}$.

Lemma 7. *Let (\mathbf{u}, λ) be a Newton-stable eigenpair of \mathcal{W} with $\lambda \geq 1$. There are $c_1, c_2, \varepsilon_0 > 0$ such that for all sufficiently small $\varepsilon > 0$ the following holds. For any $\tilde{\mathcal{W}}$ such that $\|\tilde{\mathcal{W}} - \mathcal{W}\|_F \leq \varepsilon$ there exists a unique eigenpair $(\tilde{\mathbf{u}}, \tilde{\lambda})$ of $\tilde{\mathcal{W}}$ such that*

$$\|\mathbf{u} - \tilde{\mathbf{u}}\| \leq c_1 \varepsilon \quad \text{and} \quad |\tilde{\lambda} - \lambda| \leq c_2 \varepsilon.$$

In addition, $(\tilde{\mathbf{u}}, \tilde{\lambda})$ is Newton-stable and any other eigenvector $\tilde{\mathbf{v}}$ of $\tilde{\mathcal{W}}$ satisfies $\|\tilde{\mathbf{v}} - \mathbf{u}\| \geq \varepsilon_0$.

Proof of Lemma 7. For a tensor $\mathcal{T} \in \mathbb{R}^{d \times d \times d}$ let $\mathbf{t} \in \mathbb{R}^s$ be the vector of $s = d^3$ entries $\{\mathcal{T}_{ijk}\}$. Define the function $\mathbf{Q} : \mathbb{R}^{d+s} \rightarrow \mathbb{R}^d$ by

$$\mathbf{Q}(\mathbf{v}, \mathbf{t}) = \mathcal{T}(I, \mathbf{v}, \mathbf{v}) - \mathcal{T}(\mathbf{v}, \mathbf{v}, \mathbf{v}) \cdot \mathbf{v}.$$

Note that for any $\mathbf{t} \in \mathbb{R}^s$ and $(\mathbf{v}, \beta) \in \mathbb{R}^d \times \mathbb{R}$ with $\mathbf{v} \neq \mathbf{0}$ and $\beta \neq 0$, we have that $\mathbf{Q}(\mathbf{v}, \mathbf{t}) = \mathbf{0}$ if and only if (\mathbf{v}, β) is an eigenpair of \mathbf{t} with eigenvalue $\beta = \mathcal{T}(\mathbf{v}, \mathbf{v}, \mathbf{v})$.⁴ Denote the gradients of \mathbf{Q} with respect to \mathbf{v} and \mathbf{t} by

$$\begin{aligned} A(\mathbf{v}, \mathbf{t}) &= \nabla_{\mathbf{v}} \mathbf{Q}(\mathbf{v}, \mathbf{t}) \in \mathbb{R}^{d \times d}, \\ B(\mathbf{v}, \mathbf{t}) &= \nabla_{\mathbf{t}} \mathbf{Q}(\mathbf{v}, \mathbf{t}) \in \mathbb{R}^{d \times s}. \end{aligned}$$

Let $\mathbf{w} \in \mathbb{R}^s$ be the vectorization of \mathcal{W} and let $(\mathbf{u}, \lambda) \in \mathbb{S}_{d-1} \times \mathbb{R}_+$ be a Newton-stable eigenpair of \mathbf{w} with $\lambda \geq 1$. Since \mathbf{u} is Newton-stable and $\lambda > 0$, $A(\mathbf{u}, \mathbf{w})$ is invertible. In addition, the following $(d+s) \times (d+s)$ matrix is invertible,

$$D(\mathbf{u}, \mathbf{w}) = \begin{pmatrix} A(\mathbf{u}, \mathbf{w}) & B(\mathbf{u}, \mathbf{w}) \\ \mathbf{0} & I_s \end{pmatrix}.$$

⁴This does not precisely hold when $\beta = 0$ since $\mathbf{Q}(\mathbf{v}, \mathbf{t}) = \mathbf{0}$ does not imply $\|\mathbf{v}\| = 1$ in this case, but only that \mathbf{v} is proportional to an eigenvector.

Let $\gamma_D = 1/\|D(\mathbf{u}, \mathbf{w})^{-1}\| > 0$ be the smallest singular value of $D(\mathbf{u}, \mathbf{w})$ and let $L_D < \infty$ be the Lipschitz constant of $\nabla \mathbf{Q}(\mathbf{v}, \mathbf{t}) = [A(\mathbf{v}, \mathbf{t}), B(\mathbf{v}, \mathbf{t})] \in \mathbb{R}^{d \times (d+s)}$ in a small neighborhood of (\mathbf{u}, \mathbf{w}) , namely, $\forall (\mathbf{v}, \mathbf{t}), (\tilde{\mathbf{v}}, \tilde{\mathbf{t}})$ in the neighborhood,

$$\|\nabla \mathbf{Q}(\mathbf{v}, \mathbf{t}) - \nabla \mathbf{Q}(\tilde{\mathbf{v}}, \tilde{\mathbf{t}})\| \leq L_D \|(\mathbf{v}, \mathbf{t}) - (\tilde{\mathbf{v}}, \tilde{\mathbf{t}})\|.$$

Let $B_\varepsilon(\mathbf{w}) \subset \mathbb{R}^s$ be the ball of radius ε centered at \mathbf{w} . Then by the implicit function theorem (Hubbard & Hubbard, 2015), for any $\varepsilon \leq \varepsilon_1 := \gamma_D^2/(2L_D)$, there exists a unique continuously differentiable mapping $\tilde{\mathbf{u}} : B_\varepsilon(\mathbf{w}) \rightarrow B_{2\varepsilon/\gamma_D}(\mathbf{u})$ such that $\mathbf{Q}(\tilde{\mathbf{u}}(\tilde{\mathbf{w}}), \tilde{\mathbf{w}}) = \mathbf{0}$ for all $\tilde{\mathbf{w}} \in B_\varepsilon(\mathbf{w})$. In other words, for any $\tilde{\mathbf{w}}$ such that $\|\tilde{\mathbf{w}} - \mathbf{w}\| \leq \varepsilon$, there exist a unique vector $\tilde{\mathbf{u}}$ in all $B_{2\varepsilon/\gamma_D}(\mathbf{u})$ that is an eigenpair of $\tilde{\mathbf{w}}$. Equivalently, for $\tilde{\mathcal{W}}$ such that $\|\tilde{\mathcal{W}} - \mathcal{W}\|_F \leq \varepsilon$, there exists a unique eigenvector $\tilde{\mathbf{u}}$ of $\tilde{\mathcal{W}}$ such that

$$\|\tilde{\mathbf{u}} - \mathbf{u}\| \leq 2\varepsilon/\gamma_D := c_1 \varepsilon. \quad (37)$$

The bound on $|\tilde{\lambda} - \lambda|$ readily follows from (37). Indeed, let $q : \mathbb{R}^{d+s} \rightarrow \mathbb{R}$ be $q(\mathbf{v}, \mathbf{t}) = \mathcal{T}(\mathbf{v}, \mathbf{v}, \mathbf{v})$ and let L_λ be the Lipschitz constant of q in the neighborhood of (\mathbf{u}, \mathbf{w}) . Then,

$$\begin{aligned} |\tilde{\lambda} - \lambda| &= |q(\tilde{\mathbf{u}}, \tilde{\mathbf{w}}) - q(\mathbf{u}, \mathbf{w})| \\ &\leq L_\lambda \sqrt{\|\tilde{\mathbf{u}} - \mathbf{u}\|^2 + \|\tilde{\mathbf{w}} - \mathbf{w}\|^2} \\ &\leq L_\lambda \sqrt{\frac{2}{\gamma_D} + 1} \cdot \varepsilon := c_2 \varepsilon. \end{aligned}$$

As for the Newton-stability of $\tilde{\mathbf{u}}$, let $r : \mathbb{R}^{d+s} \rightarrow \mathbb{R}_+$ be $r(\mathbf{v}, \mathbf{t}) = 1/\|A(\mathbf{v}, \mathbf{t})^{-1}\|$, the minimal singular value of $A(\mathbf{v}, \mathbf{t})$. Since (\mathbf{u}, λ) is a Newton-stable eigenpair of \mathbf{w} , $\exists \gamma_A > 0$ such that $r(\mathbf{u}, \mathbf{w}) \geq \gamma_A$. Let L_γ be the Lipschitz constant of $r(\mathbf{v}, \mathbf{t})$ in the neighborhood (Golub & Van Loan, 2012). Then, for $\varepsilon \leq \varepsilon_2 := \gamma/(2L_\gamma)$, we have $r(\tilde{\mathbf{u}}, \tilde{\mathbf{w}}) \geq \gamma_A/2 > 0$, so $(\tilde{\mathbf{u}}, \tilde{\lambda})$ is a Newton-stable eigenpair of $\tilde{\mathbf{w}}$.

Finally, we show that any other eigenvector $\tilde{\mathbf{v}}$ of $\tilde{\mathcal{W}}$ is apart from \mathbf{u} . Since $\tilde{\mathbf{u}}$ is Newton-stable, there exists $\varepsilon_0 > 0$ such that $\|\tilde{\mathbf{v}} - \tilde{\mathbf{u}}\| \geq 2\varepsilon_0$ for any other eigenvector $\tilde{\mathbf{v}}$. Hence, for $\varepsilon \leq \varepsilon_0$,

$$\|\tilde{\mathbf{v}} - \mathbf{u}\| \geq \|\tilde{\mathbf{v}} - \tilde{\mathbf{u}}\| - \|\tilde{\mathbf{u}} - \mathbf{u}\| \geq \|\tilde{\mathbf{v}} - \tilde{\mathbf{u}}\| - \varepsilon \geq \varepsilon_0.$$

Taking $\varepsilon \leq \min\{\varepsilon_0, \varepsilon_1, \varepsilon_2\}$ and c_1, c_2, ε_0 as above concludes the proof of the lemma.

Lastly, for completeness, we show that $\gamma_D \geq \frac{\gamma_A}{\sqrt{\gamma_A^2 + d}}$.

$$\begin{aligned} \gamma_D^{-1} &= \|D(\mathbf{u}, \mathbf{w})^{-1}\| \\ &\leq \sqrt{\|A(\mathbf{u}, \mathbf{w})^{-1}\|^2 (1 + \|B(\mathbf{u}, \mathbf{w})\|^2) + \|I_s\|^2} \\ &\leq \sqrt{1 + \frac{1 + \|B(\mathbf{u}, \mathbf{w})\|^2}{\gamma_A^2}}. \end{aligned} \quad (38)$$

To bound $\|B(\mathbf{u}, \mathbf{w})\|$, note that $\mathbf{Q}(\mathbf{u}, \mathbf{w})$ is linear in \mathbf{w} and its i -th entry is given by

$$[\mathbf{Q}(\mathbf{u}, \mathbf{w})]_i = \sum_{k,l} w_{ikl} u_k u_l - \left(\sum_{j,k,l} w_{jkl} u_j u_k u_l \right) u_i.$$

Thus, the $d \times m$ matrix $B(\mathbf{u}, \mathbf{w})$ has entries

$$[B(\mathbf{u}, \mathbf{w})]_{i,(jkl)} = [\nabla_{\mathbf{w}} \mathbf{Q}(\mathbf{u}, \mathbf{w})]_{i,(jkl)} = (\delta_{ij} - u_i u_j) u_k u_l,$$

which is independent of \mathbf{w} . Recalling that $\|\mathbf{u}\| = 1$,

$$\begin{aligned} \|B(\mathbf{u})\|^2 &\leq \|B(\mathbf{u})\|_F^2 = \sum_{i,j,k,l=1}^d (\delta_{ij} - u_i u_j)^2 u_k^2 u_l^2 \\ &= \sum_{i,j=1}^d (\delta_{ij}^2 - 2\delta_{ij} u_i u_j + u_i^2 u_j^2) = d - 1. \end{aligned}$$

Putting this bound in (38), we obtain $\gamma_D \geq \frac{\gamma_A}{\sqrt{\gamma_A^2 + d}}$. \square

F. Proof of Lemma 4

Let $\mathbf{v}^* \in W^\dagger$. Then $\exists i \in [d]$ such that $\mathbf{v}^{*\top} W^\top \mathbf{h} = h_i \in \{0, 1\}$. Hence, by (28), the c.d.f. $F_{\mathbf{v}^*}$ of $\mathbf{v}^{*\top} \mathbf{x}$ corresponds to the two component GMM

$$(1 - p_i) \cdot \mathcal{N}(0, \sigma^2 \|\mathbf{v}^*\|^2) + p_i \cdot \mathcal{N}(1, \sigma^2 \|\mathbf{v}^*\|^2).$$

By Lemma 1 we have $p_i = 1/\lambda(\mathbf{v}^*)^2$. Thus, $F_{\mathbf{v}^*} = G_{\mathbf{v}^*}$.

For the other direction, let $\mathbf{v} \in V_\sigma \setminus W^\dagger$. Since W is full rank, the d -dimensional vector $\mathbf{u}^\top = \mathbf{v}^\top W^\top \notin \{\mathbf{e}_i^\top\}_{i=1}^d$. Moreover, by Eq. (23) of Lemma 3,

$$\inf_{\mathbf{v} \in V_\sigma \setminus W^\dagger} \min_{\mathbf{v}^* \in W^\dagger} \|\mathbf{v} - \mathbf{v}^*\| \geq \delta_1 > 0.$$

Hence, there exists $\varepsilon_0 > 0$ such that

$$\min_{i \in [d]} \|\mathbf{u} - \mathbf{e}_i\| \geq \varepsilon_0.$$

So by the expected rigidity condition (27), there exists $\eta_0 > 0$ such that $r(\mathbf{u}) \geq \eta_0$. It follows that $F_{\mathbf{v}}$ has a component with mean that is bounded away from both 0 and 1 and thus $F_{\mathbf{v}} \neq G_{\mathbf{v}}$. In particular, there exists $\eta_1 > 0$ such that

$$\sup_{t \in \mathbb{R}} |F_{\mathbf{v}}(t) - G_{\mathbf{v}}(t)| \geq \eta_1.$$

G. Proof of Lemma 5

Recall that our sample of size $2n$ was split into two separate parts each of size n . The first n samples were used to estimate the tensor eigenvectors, and the last n samples to estimate the empirical cdf's of their projections onto the eigenvectors.

For any $\hat{\mathbf{v}}$ that is close to a vector \mathbf{v} , we bound $\Delta_n(\hat{\mathbf{v}}) = \|\hat{F}_{\hat{\mathbf{v}}} - G_{\hat{\mathbf{v}}}\|_\infty$ by the triangle inequality,

$$\begin{aligned} \|\hat{F}_{\hat{\mathbf{v}}} - G_{\hat{\mathbf{v}}}\|_\infty &\leq \|\hat{F}_{\hat{\mathbf{v}}} - F_{\hat{\mathbf{v}}}\|_\infty + \|F_{\hat{\mathbf{v}}} - F_{\mathbf{v}}\|_\infty \\ &\quad + \|F_{\mathbf{v}} - G_{\mathbf{v}}\|_\infty + \|G_{\mathbf{v}} - G_{\hat{\mathbf{v}}}\|_\infty. \end{aligned} \quad (39)$$

We now consider each of the four terms separately, starting with the first one. Since $\sigma > 0$, the cdf $F_{\hat{\mathbf{v}}} : \mathbb{R} \rightarrow [0, 1]$ is continuous and the distribution of $\|\hat{F}_{\hat{\mathbf{v}}} - F_{\hat{\mathbf{v}}}\|_\infty$ is independent of $\hat{\mathbf{v}}$. Then, by the Dvoretzky-Kiefer-Wolfowitz inequality, $\|\hat{F}_{\hat{\mathbf{v}}} - F_{\hat{\mathbf{v}}}\|_\infty$ is w.h.p. of order $O(1/\sqrt{n})$ for any $\hat{\mathbf{v}}$, and in particular tends to zero as $n \rightarrow 0$.

As for the second term, write $\hat{\mathbf{v}} = \mathbf{v} + \boldsymbol{\eta}$. Then,

$$\hat{\mathbf{v}}^\top \mathbf{x} = \mathbf{v}^\top \mathbf{x} + \boldsymbol{\eta}^\top \mathbf{x}.$$

Recall that $\mathbf{x} = W^\top \mathbf{h} + \sigma \boldsymbol{\xi}$. Hence, $|\boldsymbol{\eta}^\top \mathbf{x}| \leq \|W\|_2 \sqrt{d} \|\boldsymbol{\eta}\| + \sigma |\boldsymbol{\eta}^\top \boldsymbol{\xi}|$. The term $\boldsymbol{\eta}^\top \boldsymbol{\xi}$ is simply a zero mean Gaussian random variable with standard deviation $\sigma \|\boldsymbol{\eta}\|$. So, there exists $K_n > \sqrt{d} \|W\|_2 + \sigma n^{1/3}$ such that with probability tending to one as $n \rightarrow \infty$, for all n samples $\mathbf{x}_j \in X$, $|\boldsymbol{\eta}^\top \mathbf{x}_j| \leq K_n \|\boldsymbol{\eta}\|$. Thus, $|\hat{\mathbf{v}}^\top \mathbf{x} - \mathbf{v}^\top \mathbf{x}|$ can be bounded by $K_n \|\hat{\mathbf{v}} - \mathbf{v}\|$. This, in turn, implies that

$$\|F_{\hat{\mathbf{v}}} - F_{\mathbf{v}}\|_\infty \leq L K_n \|\hat{\mathbf{v}} - \mathbf{v}\|,$$

where $L = \max_t F'_v(t)$, which is finite for any $\sigma > 0$. Now, suppose the sequence $\hat{\mathbf{v}}_{(n)}$ converges to some \mathbf{v} at rate $O_P(1/\sqrt{n})$. Since K_n grows much more slowly with n , this term tends to zero.

Let us next consider the fourth term, and leave the third term to the end. Here note that $G_{\mathbf{v}}$ is continuous in its parameter \mathbf{v} . So if the sequence $\hat{\mathbf{v}}_{(n)}$ converges to some \mathbf{v} , then this term tends to zero.

Finally, consider the third term. If the limiting vector \mathbf{v} belongs to the correct set, namely $\mathbf{v}^* \in W^\dagger$, then $F_{\mathbf{v}} = G_{\mathbf{v}}$, and thus overall $\|\hat{F}_{\hat{\mathbf{v}}} - G_{\hat{\mathbf{v}}}\|_\infty$ tends to zero as required.

In contrast, if $\hat{\mathbf{v}}$ converges to a vector $\mathbf{v} \notin W^\dagger$, then instead of Eq. (39) we invoke the following inequality:

$$\begin{aligned} \|\hat{F}_{\hat{\mathbf{v}}} - G_{\hat{\mathbf{v}}}\|_\infty &\geq \|F_{\mathbf{v}} - G_{\mathbf{v}}\|_\infty - \|F_{\mathbf{v}} - F_{\hat{\mathbf{v}}}\|_\infty \\ &\quad - \|F_{\hat{\mathbf{v}}} - \hat{F}_{\hat{\mathbf{v}}}\|_\infty - \|G_{\hat{\mathbf{v}}} - G_{\mathbf{v}}\|_\infty. \end{aligned}$$

Here $\|F_{\mathbf{v}} - G_{\mathbf{v}}\|_\infty$ is strictly larger than zero whereas the three other remaining terms tend to zero as $n \rightarrow \infty$ as above.

H. Proof of Lemma 6

Multiplying (35) from the right by the full rank matrix Y^{-1} we obtain the equations

$$[\mathcal{W}(Y, Y, I)]_{ij} = [Y^\top]_{ij}, \quad \forall i, j \in [d].$$

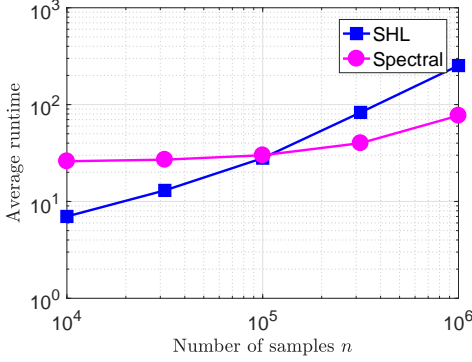


Figure 5. Average runtime vs. sample size n

Note that for all $i \in [d]$,

$$[\mathcal{W}(Y, Y, I)]_{ii} = [\mathcal{W}(\mathbf{y}_i, \mathbf{y}_i, I)]_j.$$

Since \mathcal{W} is symmetric, we thus have

$$\mathcal{W}(I, \mathbf{y}_i, \mathbf{y}_i) = \mathbf{y}_i, \quad \forall i \in [d].$$

Writing $\mathbf{y}_i = \mathbf{u}_i / \lambda_i$ we obtain the eigenpair equation

$$\mathcal{W}(I, \mathbf{u}_i, \mathbf{u}_i) = \lambda_i \mathbf{u}_i, \quad \forall i \in [d].$$

The other direction readily follows from the definition of eigenpairs.

I. Simulation runtime

Figure 5 shows the simulation runtime of the spectral approach and that of SHL vs. the number of samples n . The setup is similar to the one described in section 5. The runtime of SHL increases linearly with n , as expected. For our spectral method, for lower values of n the dominant factor is the computation of tensor eigenvectors, which does not depend on n . For large n , the dominant factor is the computation of the correlation tensor, linear in n .

J. Matrix and tensor denoising

In Algorithm 1, we modify the diagonal elements of M, \mathcal{M} by (17). This modification is suited for additive Gaussian noise, but is not applicable for the case where $X = \text{binomial}(2, W^T H)$. Instead, we implemented a method derived in (Jain & Oh, 2014) for a similar setup.

First, we treat the diagonal elements of M_σ as missing data, and complete them with the following iterative steps. (i) compute the first d eigenpairs $\{\mathbf{v}_i, \lambda_i\}$ of $R^{(k)}$; and (ii) update the diagonal elements by $R_{jj}^{(k+1)} = (\sum_i \lambda_i \mathbf{v}_i \mathbf{v}_i^T)_{jj}$.

Next, instead of computing \mathcal{M}_σ via (17) and then \mathcal{W}_σ via (19), we compute \mathcal{W}_σ directly by solving the following system of linear equations. Let K^\dagger be the pseudo-inverse

matrix of K , and $P_\Omega(\mathcal{T})$ denote a masking operation over the tensor \mathcal{T} such that,

$$P_\Omega(\mathcal{T}) = \begin{cases} \mathcal{T}_{ijk} & i \neq j \neq k \\ 0 & \text{o.w} \end{cases}$$

We estimate \mathcal{W} by the following minimization problem,

$$\hat{\mathcal{W}} = \underset{\mathcal{W}}{\operatorname{argmin}} \|P_\Omega(\mathcal{W}(K^\dagger, K^\dagger, K^\dagger)) - P_\Omega(\mathcal{M})\|_F^2$$

This method depends only on the off-diagonal elements of M and \mathcal{M} and hence is applicable whenever $\mathbb{E}[X|H] = W^T H$ and the noise has bounded variance.

K. Adding a weighted least square step to the spectral method

In section 5, we compare the results of algorithm 1 with and without an additional single weighted least square step. Given an estimate \hat{W} , for each observed instance \mathbf{x}_j we calculate the conditional likelihood $\mathcal{L}(\mathbf{x}_j | \mathbf{h})$ for the 2^d possible binary vectors $\mathbf{h} \in \{0, 1\}^d$,

$$\mathcal{L}(\mathbf{x}_j | \mathbf{h}) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\|\mathbf{x}_j - \hat{W}^T \mathbf{h}\|^2 / (2\sigma^2)\right)$$

For each instance \mathbf{x}_j , we keep the top $K = 6$ vectors $\mathbf{h}_{1j}, \dots, \mathbf{h}_{Kj}$ with the highest likelihood. Let $\Pi \in [0, 1]^{K \times n}$ be a weight matrix such that Π_{kj} is proportional to $\mathcal{L}(\mathbf{x}_j | \mathbf{h}_{kj})$, and $\sum_k \Pi_{kj} = 1$ for all j . The new estimate \hat{W}_{wls} is the minimizer of the weighted least square problem,

$$\hat{W}_{\text{wls}} = \underset{W}{\operatorname{argmin}} \sum_{j=1}^n \sum_{k=1}^K \Pi_{kj} \|\mathbf{x}_j - W^T \mathbf{h}_{kj}\|^2.$$

L. Alternating least squares for W and H

In section 5, we compare the results of the spectral approach to the following ALS iterations, with a random starting point.

$$\begin{aligned} W^{(k)} &= \underset{W \in \mathbb{R}^{d \times m}}{\operatorname{argmin}} \|X - W^T H^{(k-1)}\|_F^2 \\ \hat{H}^{(k)} &= \underset{H \in \mathbb{R}^{d \times n}}{\operatorname{argmin}} \|X - (W^{(k)})^T H\|_F^2 \\ H^{(k)} &= \underset{H \in \{0,1\}^{d \times n}}{\operatorname{argmin}} \|H - \hat{H}^{(k)}\|_F^2, \end{aligned}$$

M. Genetic admixture simulations

The simulated admixture data was generated via the following steps:

1. We used SCRM (Staab et al., 2015) to simulate a split between $d = 3$ ancestral populations, with separation

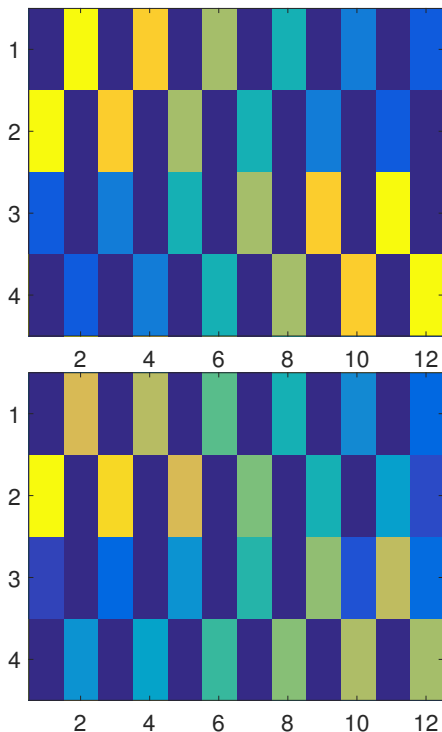


Figure 6. Real (upper panel) and estimated (lower panel) cell proportion matrix

time of 4000 generations. The simulator generated 40 chromosomes of length $250 \cdot 10^6$ for each of the three populations. The simulation parameters were determined as $N_0 = 10^4$ effective population size, 10^{-8} mutation rate (per base pair per generation), and 10^{-8} recombination rate (per base pair per generation).

2. We sampled the proportion matrix W from a Dirichlet distribution with parameter α .
3. Two chromosomes of length $250 \cdot 10^6$ were created for each of the $m = 50$ admixed individuals with the following steps: (i) An ancestral population was sampled according to W , say, population h_A . (ii) One of the 40 chromosomes was sampled from h_A , say $h_A(k)$ (iii) A block length l was sampled from an exponential distribution with rate 20 per Morgan corresponding admixture event happening 20 generations ago (in our case, 1 Morgan was 10^8 base pairs). (iv) A block of length l was copied from chromosome $h_A(k)$ to the corresponding locations in the new admixed chromosome. We repeated steps (i)-(iv) until completion of the chromosome.

N. Analysis of DNA methylation data

The dataset is part of the supplementary material of (Houseman et al., 2012). The observed matrix $X \in$

$[0, 1]^{m \times n}$ consists of $m = 12$ blood samples, each with the DNA methylation measurements in $n = 500$ sites (called CpGs).

The statistical model for X is similar to that of Eq. (1). We assume that each blood cell is a mixture of $d = 4$ cell types, with unknown proportions. The latent variables h correspond to the presence or absence of methylation in each site for the 4 cell types. Given the DNA methylation array, the task is to estimate the proportion matrix W .

The upper and lower panels of Figure 6 correspond to the real and estimated mixture matrix. For comparison, we performed the steps described in detail for this dataset in Slawski et al. (2013, Section 4). For both methods we measured the l_1 distance between the real and estimated mixture matrices,

$$\frac{1}{mn} \sum_{ij} |W_{ij} - \hat{W}_{ij}|$$

The l_1 distance were equal to 0.003 for SHL and 0.00056 for the spectral approach.