# Maximum Likelihood for Matrices with Rank Constraints

Jonathan Hauenstein[1], Jose Israel Rodriguez[2], Bernd Sturmfels [2]

[1] *Department of Mathematics, North Carolina State University, Raleigh, NC 27695, USA*
[2] *Department of Mathematics, University of California, Berkeley, CA 94720, USA,*

**Abstract.** Maximum likelihood estimation is a fundamental optimization problem in statistics. We study this problem on manifolds of matrices with bounded rank. These represent mixtures of distributions of two independent discrete random variables. We determine the maximum likelihood degree for a range of determinantal varieties, and we apply numerical algebraic geometry to compute all critical points of their likelihood functions. This led to the discovery of maximum likelihood duality between matrices of complementary ranks, a result proved subsequently by Draisma and Rodriguez.
**2000 Mathematics Subject Classifications**: Maximum likelihood degree, rank constraint, determinantal variety
**Key Words and Phrases**: 62F10, 13P15, 14M12, 90C26

## 1. Introduction

Maximum likelihood estimation (MLE) is a fundamental computational task in statistics. A typical problem encountered in its applications is the occurrence of multiple local maxima. In order to be certain that a global maximum of the likelihood function has been achieved, one needs to locate all solutions to a system of polynomial equations. In this paper we study these equations for two discrete random variables, having $m$ and $n$ states respectively. A joint probability distribution for two such random variables is written as an $m \times n$-matrix:

$$P \quad = \quad \begin{pmatrix} p_{11} & p_{12} & \cdots & p_{1n} \\ p_{21} & p_{22} & \cdots & p_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ p_{m1} & p_{m2} & \cdots & p_{mn} \end{pmatrix}. \tag{1}$$

The entry $p_{ij}$ represents the probability that the first variable is in state $i$ and the second is in state $j$. Thus, the entries of $P$ are non-negative and their sum $p_{++}$ is 1. By a statistical model, we mean a closed subset $\mathcal{M}$ of the probability simplex $\Delta_{mn-1}$ of all such matrices $P$.

If i.i.d. samples are drawn from some $P$ then we summarize the data also in a matrix

$$U \quad = \quad \begin{pmatrix} u_{11} & u_{12} & \cdots & u_{1n} \\ u_{21} & u_{22} & \cdots & u_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ u_{m1} & u_{m2} & \cdots & u_{mn} \end{pmatrix}. \tag{2}$$

The entries of $U$ are non-negative integers whose sum is $u_{++}$. As is customary in algebraic statistics [9, 15, 26], we write the *likelihood function* corresponding to the data matrix $U$ as

$$\ell_U \quad = \quad \frac{\prod_{i=1}^{m} \prod_{j=1}^{n} p_{ij}^{u_{ij}}}{\left( \sum_{i=1}^{m} \sum_{j=1}^{n} p_{ij} \right)^{u_{++}}}. \tag{3}$$

*Email addresses:* hauenstein@ncsu.edu (J. Hauenstein), jo.ro@berkeley.edu (J. I. Rodriguez), bernd@math.berkeley.edu (J.Sturmfels)

This formula defines a rational function on the complex projective space $\mathbb{P}^{mn-1}$ whose restriction to the simplex $\Delta_{mn-1}$ is the usual likelihood function divided by a multinomial coefficient. The MLE problem is to find the global maximum of $\ell_U$ over the model $\mathcal{M}$.

Our model of interest is the set $\mathcal{M}_r$ of matrices $P$ of rank $\leq r$. This is the intersection of the variety $\mathcal{V}_r \subset \mathbb{P}^{mn-1}$ defined by the $(r+1) \times (r+1)$-minors of $P$ with $\Delta_{mn-1}$. For generic $U$, the rational function $\ell_U$ has finitely many critical points on the determinantal variety $\mathcal{V}_r$. Their number is the *ML degree* of $\mathcal{V}_r$. In this paper, we formulate a polynomial system consisting of $mn$ equations in $mn$ variables defining such critical points and compute them using methods from numerical algebraic geometry. That computation enables us to reliably find all local maxima of the likelihood function $\ell_U$ among positive points in $\mathcal{M}_r$. Among the new results is the determination of the bold face numbers in the following table.

**Theorem 1.** *The known values for the ML degrees of the determinantal varieties $\mathcal{V}_r$ are*

| $(m,n) =$ | $(3,3)$ | $(3,4)$ | $(3,5)$ | $(4,4)$ | $(4,5)$ | $(4,6)$ | $(5,5)$ | |
|---|---|---|---|---|---|---|---|---|
| $r = 1$ | 1 | 1 | 1 | 1 | 1 | 1 | 1 | |
| $r = 2$ | 10 | 26 | **58** | **191** | **843** | **3119** | **6776** | (4) |
| $r = 3$ | 1 | 1 | 1 | **191** | **843** | **3119** | **61326** | |
| $r = 4$ | | | | 1 | 1 | 1 | **6776** | |
| $r = 5$ | | | | | | | 1 | |

The smaller numbers 10 and 26 had already been computed in [15, §5], but the symbolic computations using `Singular` that were presented in [15] had failed beyond the size $3 \times 4$.

In 2005, the third author offered a cash prize of 100 Swiss Francs (cf. [26, §3]) for the solution of a particular $4 \times 4$-instance that was described in [21, Example 1.16]. That prize was won in 2008 by Mingfu Zhu who solved this challenge in [29]. See also [24, Example 5.2] for a solution using `Singular`, and [10] for a statistical perspective on this problem. However, none of these papers had found the number 191 of critical points for the $4 \times 4$ cases. In the first version of this paper, we stated the conjecture that the column symmetry among the ML degrees always holds. This has subsequently been proven by Draisma and Rodriguez:

**Theorem 2** ([8]). *If $m \leq n$ then the ML degrees for rank $r$ and for rank $m - r + 1$ coincide.*

Our findings might appeal also to those interested in the topology of algebraic varieties. For a variety $\mathcal{V}$ in $\mathbb{P}^{mn-1}$, let $\mathcal{V}^0$ denote the open subset given by $p_{11}p_{12} \cdots p_{mn}p_{++} \neq 0$. Huh [16] recently proved that if $\mathcal{V}^0$ is smooth then the ML degree of $\mathcal{V}$ is equal to the signed Euler characteristic of $\mathcal{V}^0$. In our case, for $r \geq 2$, the open determinantal variety $\mathcal{V}_r^0$ is singular along $\mathcal{V}_{r-1}^0$, but a suitably modified statement is expected to be true. It might be speculated that the results in Theorems 1 and 2 will ultimately have a topological explanation.

The entries "1" of the table in (4) have easy explanations. For $r = m$ we have $\mathcal{V}_m = \mathbb{P}^{mn-1}$ and the unique critical point of the likelihood function $\ell_U$ is $P = \frac{1}{u_{++}}U$. The first row of (4) states that the independence model $\mathcal{M}_1$ has ML degree 1. This fact is well-known to statisticians, as the rank 1 matrix with entries $(u_{i+}u_{+j})/u_{++}^2$ is the unique critical point for $\ell_U$ on $\mathcal{V}_1^0$. We found it instructive to derive this fact from Huh's result [16, Theorem 1.(iii)]:

**Example 1.** Let $r = 1$. The Segre variety $\mathcal{V}_1 = \mathbb{P}^{m-1} \times \mathbb{P}^{n-1}$ is smooth. Fix coordinates $(x_1 : \cdots : x_m)$ on $\mathbb{P}^{m-1}$ and coordinates $(y_1 : \cdots : y_n)$ on $\mathbb{P}^{n-1}$. The open subset $\mathcal{V}_1^0$ consists of all points in $\mathbb{P}^{m-1} \times \mathbb{P}^{n-1}$ with $x_1 x_2 \cdots x_m y_1 y_2 \cdots y_n (x_1 + \cdots + x_m)(y_1 + \cdots + y_n) \neq 0$. Hence

$$\mathcal{V}_1^0 = (\mathbb{P}^{m-1} \text{ minus } m + 1 \text{ hyperplanes}) \times (\mathbb{P}^{n-1} \text{ minus } n + 1 \text{ hyperplanes}).$$

Each factor has signed Euler characteristic 1, and hence so does their product. $\square$

This article is organized as follows. In Section 2, we formulate the constraints that characterize critical points of $\ell_U$ on $\mathcal{V}_r$ as a square system of polynomial equations. The specific

formulation in Theorem 3 is one of our key contributions. It is used to derive upper bounds in terms of $m$, $n$, and $r$. Theorem 4 extends our results to the case of symmetric matrices, and hence to mixtures of two identically distributed random variables.

Section 3 is devoted to our computations using numerical algebraic geometry. This furnishes valuable new tools for practitioners of statistics who are interested in exploring probability one algorithms for computing the global maximum of a given likelihood function.

In Section 4, we introduce a refined version of Theorem 2, now also proved in [8], and we summarize the computational evidence we had gathered to support it. The Galois group computations in Proposition 1 might be of independent interest. In Theorem 8, we present a proof of [29, Conjecture 11] by means of certified numerical computations.

Section 5 features the statistical view on our approach, and we explain how it differs from running the EM algorithm for discrete mixture models. The determinantal variety $\mathcal{V}_r$ is the Zariski closure of the latent variable model for $r$-fold mixtures of independent variables. They are equal in $\Delta_{mn-1}$ if and only if $r \leq 2$. For $r \geq 3$ this takes us to the real algebraic geometry problem, pioneered in [19], of distinguishing between rank and non-negative rank.

## 2. Equations and bounds

In this section, we present several formulations of the critical equations for the likelihood function on the determinantal variety $\mathcal{V}_r = \{\mathrm{rank}(P) \leq r\}$. We view $\mathcal{V}_r$ as an affine variety in the space of matrices $\mathbb{C}^{m \times n}$ and we assume $m \leq n$. Our main result is Theorem 3 which expresses our problem as a square system of $mn$ polynomial equations in $mn$ unknowns.

An $m \times n$-matrix $P$ is a regular point in the determinantal variety $\mathcal{V}_r$ if and only if $\mathrm{rank}(P) = r$. If this holds then the tangent space $T_P$ is a linear subspace of dimension $rn + rm - r^2$ in $\mathbb{C}^{m \times n}$, and its orthogonal complement (with respect to the standard inner product) is a linear subspace $T_P^\perp$ of dimension $(m - r)(n - r)$ in $\mathbb{C}^{m \times n}$.

Our input is a strictly positive data matrix $U$. We consider the logarithm of the likelihood function $\ell_U$ as in (3). The partial derivatives of the *log-likelihood function* $\log(\ell_U)$ are then

$$\frac{\partial \log(\ell_U)}{\partial p_{ij}} \;=\; \frac{u_{ij}}{p_{ij}} - \frac{u_{++}}{p_{++}}. \tag{5}$$

By [15, Proposition 3], a matrix $P$ of rank $r$ is a critical point for $\log(\ell_U)$ on $\mathcal{V}_r$ if and only if the linear subspace $T_P^\perp$ contains the $m \times n$-matrix whose $(i,j)$ entry is (5). Hence the system of equations we seek to solve can be expressed in the following *geometric formulation*:

$$\mathrm{rank}(P) = r\,, \qquad p_{++} = 1\,, \quad \text{and} \quad \text{the matrix } \big(u_{ij}/p_{ij} - u_{++}\big) \text{ lies in } T_P^\perp. \tag{6}$$

This is saying that the gradient of the objective function must be orthogonal to the tangent space of the variety at a critical point as in the elementary Lagrange multipliers method. When translating (6) into polynomial equations, we need to make sure to exclude matrices $P$ of rank strictly less than $r$, as these are singular points in $\mathcal{V}_r$. We also need to exclude matrices $P$ with $p_{ij} = 0$ for some $(i,j)$. These non-degeneracy conditions require some care.

In [15], the following formulation was used to represent our problem. Let $J(P)$ denote the Jacobian matrix of the prime ideal defining $\mathcal{V}_r$. Since that ideal is minimally generated by the $\binom{m}{r+1}\binom{n}{r+1}$ subdeterminants of format $(r+1) \times (r+1)$, the Jacobian $J(P)$ is a matrix of format $\binom{m}{r+1}\binom{n}{r+1} \times mn$ whose entries are homogeneous polynomials of degree $r$. Let $[U]$ denote the matrix $U$ when written as a row vector of format $1 \times mn$, and similarly $[P]$ is the vectorization of $P$. We write $\mathrm{diag}[P]$ for the diagonal $mn \times mn$-matrix with entries $p_{11}, p_{12}, \ldots, p_{mn}$. The following extended Jacobian has $2 + \binom{m}{r+1}\binom{n}{r+1}$ rows and $mn$ columns:

$$\mathcal{J}(P) \;=\; \begin{pmatrix} [U] \\ [P] \\ J(P) \cdot \mathrm{diag}[P] \end{pmatrix}.$$

For a matrix $P$ of rank $r$, the Jacobian $J(P)$ has rank $(m-r)(n-r) = \mathrm{codim}(\mathcal{V}_r)$. The third condition in (6) translates into the requirement that the span of the first two rows intersects the rowspace of $J(P) \cdot \mathrm{diag}[P]$. From this we derive the *rank formulation*

$$\mathrm{rank}(P) \le r \quad \text{and} \quad \mathrm{rank}(\mathcal{J}(P)) \le (m-r)(n-r)+1. \tag{7}$$

This formulation of our problem is elegant and is adapted to projective geometry in $\mathbb{P}^{mn-1}$. In terms of equations, we simply take the minors of size $r+1$ of the matrix $P$, and the minors of size $(m-r)(n-r)+2$ of the matrix $\mathcal{J}(P)$. However, this has two serious disadvantages: first, the number of minors is enormous, and second, we must get rid of extraneous solutions by saturation. Namely, to get rid of solutions $P$ with $\mathrm{rank}(P) \le r-1$, we need to saturate by the $r \times r$-minors of $P$, and to get rid of solutions on the boundary, we need to saturate by the product of linear forms $p_{11}p_{12}\cdots p_{mn}p_{++}$. This was done symbolically in [15, §4].

The calculation can be sped up a little bit by taking only $(m-r)(n-r)$ of the rows of $J(P)$, while also imposing the non-homogeneous equation $p_{++} = 1$. Finally, we can replace the first two rows of $J(P)$ by a single row $[U] - u_{++}[P]$ and require that the maximal minors of the resulting $((m-r)(n-r)+1) \times mn$-matrix be zero. This leads to some improvements but is still far from sufficient to get to the full range of ML degrees reported in Theorem 1.

To get to those results, we pursue the following alternatives: first, we introduce new unknowns which allow us to replace the rank conditions by bilinear equations, and, second, we represent the subspace $T_P^\perp = \mathrm{rowspace}(J(P))$ using those same new unknowns. Let $L$ be an $(m-r) \times m$-matrix of unknowns, let $R$ be an $n \times (n-r)$-matrix of unknowns, and $\Lambda = (\lambda_{ij})$ an $(n-r) \times (m-r)$-matrix of unknowns. Then our *general kernel formulation* is:

$$p_{++} = 1, \quad L \cdot P = 0, \quad P \cdot R = 0, \quad \text{and} \quad P \star (R \cdot \Lambda \cdot L)^T + u_{++} \cdot P = U. \tag{8}$$

Here $A \star B$ denotes the Hadamard (entry-wise) product of two matrices of the same format. If the rows of $L$ are linearly independent and the columns of $R$ are linearly independent, then either of the conditions $L \cdot P = 0$ and $P \cdot R = 0$ suffice to imply that $\mathrm{rank}(P) \le r$.

We now explain the last condition in (8). The space $T_P^\perp$ is spanned by the rank 1 matrices $(\rho_i \cdot \ell_j)^T$ where $\rho_i$ is the $i$-th column of $R$ and $\ell_j$ is the $j$-th row of $L$. Then

$$(R \cdot \Lambda \cdot L)^T = \sum_{i=1}^{n-r}\sum_{j=1}^{m-r} \lambda_{ij}(\rho_i \cdot \ell_j)^T$$

is a general matrix in $T_P^\perp$. The matrix $(u_{ij}/p_{ij} - u_{++})$ in (6) can be written as

$$P^{\star(-1)} \star U - u_{++} \cdot \mathbf{1}. \tag{9}$$

Hence the last condition of (6) is equivalent to saying (9) equals $(R \cdot \Lambda \cdot L)^T$ for some $\Lambda$. We write this as $(R \cdot \Lambda \cdot L)^T + u_{++} \cdot \mathbf{1} = P^{\star(-1)} \star U$. We take Hadamard product of both sides with the matrix $P$ to get the last equation in (8). This operation is invertible since all entries of $U$ are non-zero. Indeed, that last equation is $P \star ((R \cdot \Lambda \cdot L)^T + u_{++} \cdot \mathbf{1}) = U$, and if this holds then all $mn$ entries of the matrix $P$ must be non-zero.

We conclude that (8) is a correct formulation of our problem provided we can ensure

$$\mathrm{rank}(L) = m-r, \quad \mathrm{rank}(R) = n-r, \quad \text{and} \quad \mathrm{rank}(P) = r.$$

We note that (8) is highly redundant as far as the number of variables is concerned. There are several ways to reduce that number. For instance, we can simply set $\lambda_{ij} = 1$ for all $i,j$. In addition, we can either replace $L$ by a single row or replace $R$ by a single column. Even after these simplifications, the critical points of $\ell_U$ on $\mathcal{V}_r$ are still represented faithfully.

After some experimentation, we found that the following simplification steps lead to the best computational results. Recall that $m \le n$. Let $P_1$ be an $r \times r$-matrix of unknowns, let $R_1$

be an $r \times (n - r)$-matrix of unknowns, and let $L_1$ be an $(m - r) \times r$-matrix of unknowns. The matrix $\Lambda = (\lambda_{ij})$ is as before. Using this notation, we take (8) with

$$L = \begin{pmatrix} L_1 & -I_{m-r} \end{pmatrix}, \quad P = \begin{pmatrix} P_1 & P_1 R_1 \\ L_1 P_1 & L_1 P_1 R_1 \end{pmatrix}, \quad \text{and} \quad R = \begin{pmatrix} R_1 \\ -I_{n-r} \end{pmatrix}, \tag{10}$$

where $I_{m-r}$ and $I_{n-r}$ are identity matrices. We call (8) with (10) the *local kernel formulation* of our problem. Note that the constraints $L \cdot P = 0$, $P \cdot R = 0$, $\text{rank}(L) = m - r$, and $\text{rank}(R) = n - r$ are automatically satisfied in this formulation. The condition $\text{rank}(P) = r$ is also implied for every solution provided $U$ is generic. Finally, the equation $p_{++} = 1$ can be removed from (8) in this formulation since $p_{++} = 1$ is equivalent to the sum of all $mn$ equations given by $P \star (R \cdot \Lambda \cdot L)^T + u_{++} \cdot P = U$. By counting equations and unknowns, we now see that our system is a square system consisting of $mn$ equations in $mn$ unknowns.

**Theorem 3.** *Let $U$ be a generic $m \times n$ data matrix with $m \leq n$. The polynomial system*

$$P \star (R \cdot \Lambda \cdot L)^T + u_{++} \cdot P = U \tag{11}$$

*consists of $mn$ equations in $mn$ unknowns given by (10). It has finitely many complex solutions $(P_1, L_1, R_1, \Lambda)$, and the corresponding $m \times n$-matrices $P$ defined by (10) are precisely the critical points of the likelihood function $\ell_U$ on the determinantal variety $\mathcal{V}_r$.*

Since the column sums of $P \star (R \cdot \Lambda \cdot L)^T$ are zero, we can further simplify the $n$ equations. For the first $m$ columns, we replace each entry on the diagonal with the column sum. For the last $n - m$ columns, we replace the last entry in the column with the column sum.

**Example 2.** To illustrate the local kernel formulation (11), we consider $m = n = 3$ with the two subcases $r = 1$ and $r = 2$. Both have nine equations in nine unknowns.

*Subcase $r = 1$:* The nine unknowns are the entries in the matrices

$$L_1 = \begin{pmatrix} l_{11} \\ l_{21} \end{pmatrix}, \quad P_1 = \begin{pmatrix} p_{11} \end{pmatrix}, \quad R_1 = \begin{pmatrix} r_{11} & r_{12} \end{pmatrix}, \quad \Lambda = \begin{pmatrix} \lambda_{11} & \lambda_{12} \\ \lambda_{21} & \lambda_{22} \end{pmatrix},$$

and the nine equations from (11) take the form

$$
\begin{aligned}
p_{11}(1 + l_{11} + l_{21}) &= (u_{11} + u_{21} + u_{31})/u_{++} \\
p_{11}r_{11}(u_{++} - l_{11}\lambda_{11} - l_{21}\lambda_{12}) &= u_{12} \\
p_{11}r_{12}(u_{++} - l_{11}\lambda_{21} - l_{21}\lambda_{22}) &= u_{13} \\
p_{11}l_{11}(u_{++} - r_{11}\lambda_{11} - r_{12}\lambda_{21}) &= u_{21} \\
p_{11}r_{11}(1 + l_{11} + l_{21}) &= (u_{12} + u_{22} + u_{32})/u_{++} \\
p_{11}l_{11}r_{12}(\lambda_{21} + u_{++}) &= u_{23} \\
p_{11}l_{21}(u_{++} - r_{11}\lambda_{12} + r_{12}\lambda_{22}) &= u_{31} \\
p_{11}l_{21}r_{11}(\lambda_{12} + u_{++}) &= u_{32} \\
p_{11}r_{12}(1 + l_{11} + l_{21}) &= (u_{13} + u_{23} + u_{33})/u_{++}.
\end{aligned}
$$

This system has a unique solution which writes the unknowns as rational functions in the $u_{ij}$.

*Subcase $r = 2$:* The nine unknowns are the entries in the matrices

$$L_1 = \begin{pmatrix} l_{11} & l_{12} \end{pmatrix}, \quad P_1 = \begin{pmatrix} p_{11} & p_{12} \\ p_{21} & p_{22} \end{pmatrix}, \quad R_1 = \begin{pmatrix} r_{11} \\ r_{21} \end{pmatrix}, \quad \Lambda = \begin{pmatrix} \lambda_{11} \end{pmatrix},$$

| $(m,n,r)$ | $(3,3,1)$ | $(3,3,2)$ | $(3,4,1)$ | $(3,4,2)$ | $(3,5,1)$ | $(3,5,2)$ |
|---|---|---|---|---|---|---|
| Bézout | 73728 | 49152 | 3538944 | 2359296 | 169869312 | 113246208 |
| 4-hom | 270 | 1350 | 840 | 29400 | 2025 | 378000 |
| linear product | 172 | 1018 | 374 | 20844 | 650 | 68586 |
| polyhedral | 6 | 53 | 10 | 472 | 15 | 2724 |
| ML Degree | 1 | 10 | 1 | 26 | 1 | 58 |

| $(m,n,r)$ | $(4,4,1)$ | $(4,4,2)$ | $(4,4,3)$ | $(4,5,1)$ | $(4,5,2)$ | $(4,5,3)$ |
|---|---|---|---|---|---|---|
| Bézout | 905969664 | 603979776 | 402653184 | 173946175488 | 115964116992 | 77309411328 |
| 4-hom | 17600 | 7276500 | 580800 | 63700 | 323723400 | 115615500 |
| linear product | 5690 | 4791168 | 224598 | 13560 | 165869606 | 58335270 |
| polyhedral | 20 | 15280 | 2847 | 35 | 241218 | 145273 |
| ML Degree | 1 | 191 | 191 | 1 | 843 | 843 |

Table 1: Comparison of upper bounds for selected $(m, n, r)$

and the nine equations take the form

$$
\begin{aligned}
p_{11}(1 + l_{11}) + p_{21}(1 + l_{12}) &= (u_{11} + u_{21} + u_{31})/u_{++} \\
p_{12}(l_{11}r_{21}\lambda_{11} + u_{++}) &= u_{12} \\
(p_{11}r_{11} + p_{12}r_{21})(u_{++} - l_{11}\lambda_{11}) &= u_{13} \\
p_{21}(l_{12}r_{11}\lambda_{11} + u_{++}) &= u_{21} \\
p_{12}(1 + l_{11}) + p_{22}(1 + l_{12}) &= (u_{12} + u_{22} + u_{32})/u_{++} \\
(p_{21}r_{11} + p_{22}r_{21})(u_{++} - l_{12}\lambda_{11}) &= u_{23} \\
(p_{11}l_{11} + p_{21}l_{12})(u_{++} - r_{11}\lambda_{11}) &= u_{31} \\
(p_{12}l_{11} + p_{22}l_{12})(u_{++} - r_{21}\lambda_{11}) &= u_{32} \\
(p_{11}r_{11} + p_{12}r_{21})(1 + l_{11}) + (p_{21}r_{11} + p_{22}r_{21})(1 + l_{12}) &= (u_{13} + u_{23} + u_{33})/u_{++}.
\end{aligned}
$$

This system has ten complex solutions for a generic data matrix $U$. In other words, the 9 unknowns $l_{..}, p_{..}, r_{..}$ and $\lambda_{11}$ are algebraic functions of degree 10 in $u_{11}, u_{12}, \ldots, u_{33}$.  □

Upper bounds on the ML degree of $\mathcal{V}$ arise from our formulation. The Bézout bound is

$$ 2^r \cdot 3^{n-r} \cdot 4^{n(m-1)}. $$

If we consider $(P_1, L_1, R_1, \Lambda)$ in the product space $\mathbb{C}^{r^2} \times \mathbb{C}^{r(m-r)} \times \mathbb{C}^{r(n-r)} \times \mathbb{C}^{(n-r)(m-r)}$, our system consists of $r$ equations of degree $(1,1,0,0)$, $n-r$ equations of degree $(1,1,1,0)$, and $n(m-1)$ equations of degree $(1,1,1,1)$. The associated 4-homogeneous Bézout bound is the coefficient of the monomial $w^{r^2} \cdot x^{r(m-r)} \cdot y^{r(n-r)} \cdot z^{(n-r)(m-r)}$ in the expression

$$ (w + x)^r \cdot (w + x + y)^{n-r} \cdot (w + x + y + z)^{n(m-1)}. $$

A refinement of the 4-homogeneous bound using the fact that each polynomial only depends upon a subset of the variables yields a *linear product bound* [28]. Finally, the *polyhedral root count* exploits the sparsity of the monomials in our system. We computed the polyhedral bound for various cases using `MixedVol` [11] in `PHC` [27]. All of the aforementioned bounds are presented in Table 1 for selected values of $m$, $n$, and $r$. When solving a polynomial system using homotopies built from these bounds, one must balance the added computational cost required for the tighter bound with the computational savings arising from that bound.

We close this section by discussing rank constraints on symmetric matrices of the form

$$
P = \begin{pmatrix}
2p_{11} & p_{12} & p_{13} & \cdots & p_{1n} \\
p_{12} & 2p_{22} & p_{23} & \cdots & p_{2n} \\
p_{13} & p_{23} & 2p_{33} & \cdots & p_{3n} \\
\vdots & \vdots & \vdots & \ddots & \vdots \\
p_{1n} & p_{2n} & p_{3n} & \cdots & 2p_{nn}
\end{pmatrix}.
\tag{12}
$$

The case $n = 3$ was treated in [15, Example 12] where its ML degree was found to be 6. It is essential that the unknowns $p_{ii}$ on the diagonal are multiplied by 2 before imposing the rank

constraints. The matrices (12) of rank one form a Veronese variety in $\mathbb{P}^{(n+2)(n-1)/2}$. This variety has ML degree 1 and represents the independence model for two identically distributed random variables on $n$ states. The case $n = 2$ is the Hardy-Weinberg curve [21, Figure 3.1]. Larger ranks $r$ correspond to the secant varieties of this Veronese variety.

**Theorem 4.** *The known values for the ML degrees of rank $r$ symmetric matrices (12) are*

$$
\begin{array}{ccccc}
n = & 3 & 4 & 5 & 6 \\
r = 1 & 1 & 1 & 1 & 1 \\
r = 2 & 6 & \mathbf{37} & 270 & \mathbf{2341} \\
r = 3 & 1 & \mathbf{37} & 1394 & \\
r = 4 & & 1 & \mathbf{270} & \\
r = 5 & & & 1 & \mathbf{2341}
\end{array}
\tag{13}
$$

Our input is a strictly positive symmetric $n \times n$-matrix $U$. The likelihood function equals

$$
\ell_U \;\; = \;\; \frac{\prod_{i \leq j} p_{ij}^{u_{ij}}}{\left( \sum_{i \leq j} p_{ij} \right)^{\sum_{i \leq j} u_{ij}}}.
\tag{14}
$$

In the statistical context, when the sum of the $p_{ij}$ entries equals 1, we have

$$
\frac{\partial \log(\ell_U)}{\partial p_{ij}} \;\; = \;\; \frac{u_{ij}}{p_{ij}} - \sum_{i \leq j} u_{ij}.
\tag{15}
$$

We compute the critical points on the variety of rank $r$ matrices (12) by adapting the formulation in Theorem 3. Let $P_1$ be a symmetric $r \times r$-matrix of unknowns where the diagonal entries are multiplied by 2 similar to (12), let $L_1$ be an $(n - r) \times r$-matrix of unknowns, and $\Lambda$ be a symmetric $(n - r) \times (n - r)$-matrix. Following (10), we define

$$
L = \begin{pmatrix} L_1 & -I_{m-r} \end{pmatrix} \quad \text{and} \quad P = \begin{pmatrix} P_1 & P_1 L_1^T \\ L_1 P_1 & L_1 P_1 L_1^T \end{pmatrix}.
\tag{16}
$$

To account for the $p_{ii}$'s not being multiplied by 2 in the likelihood function, let $D$ be the $n \times n$-matrix whose diagonal entries are 2 and off-diagonal entries are 1. The *symmetric local kernel formulation* is the square system consisting of the upper triangular part of

$$
P \star (L^T \cdot \Lambda \cdot L) + \sum_{i \leq j} u_{ij} \cdot P \;\; = \;\; D \star U.
\tag{17}
$$

This is a system of $n(n + 1)/2$ equations in $n(n + 1)/2$ unknowns. Similar to the local kernel formulation, the column sums of $P \star (L^T \cdot \Lambda \cdot L)$ are zero. Hence (17) implies $\sum_{i \leq j} p_{ij} = 1$. We use this fact to replace the diagonal entries in (17) with the corresponding column sum.

**Example 3.** We illustrate the symmetric local kernel formulation (17) for the two subcases $r = 1, 2$ when $n = 3$. Both have 6 equations in 6 unknowns. Here, $u_{++} = \sum_{i \leq j} u_{ij}$.

*Subcase $r = 1$:* The six unknowns arise from the entries in the matrices

$$
L_1 = \begin{pmatrix} l_{11} \\ l_{21} \end{pmatrix}, \quad P_1 = \begin{pmatrix} 2p_{11} \end{pmatrix}, \quad \Lambda = \begin{pmatrix} \lambda_{11} & \lambda_{12} \\ \lambda_{12} & \lambda_{22} \end{pmatrix},
$$

and the six equations take the form

$$
\begin{aligned}
2p_{11}(1 + l_{11} + l_{21}) &= (2u_{11} + u_{12} + u_{13})/u_{++} \\
2p_{11}l_{11}(u_{++} - l_{11}\lambda_{11} - l_{21}\lambda_{12}) &= u_{12} \\
2p_{11}l_{21}(u_{++} - l_{11}\lambda_{12} - l_{21}\lambda_{22}) &= u_{13} \\
2p_{11}l_{11}(1 + l_{11} + l_{21}) &= (u_{12} + 2u_{22} + u_{23})/u_{++} \\
2p_{11}l_{11}l_{21}(\lambda_{12} + u_{++}) &= u_{23} \\
2p_{11}l_{21}(1 + l_{11} + l_{21}) &= (u_{13} + u_{23} + 2u_{33})/u_{++}.
\end{aligned}
$$

This system has a unique solution which writes the unknowns as rational functions in the $u_{ij}$.

*Subcase $r = 2$*: The six unknowns arise from the entries in the matrices

$$L_1 = \begin{pmatrix} l_{11} & l_{12} \end{pmatrix}, \quad P_1 = \begin{pmatrix} 2p_{11} & p_{12} \\ p_{12} & 2p_{22} \end{pmatrix}, \quad \Lambda = \begin{pmatrix} \lambda_{11} \end{pmatrix},$$

and the six equations take the form

$$
\begin{aligned}
2p_{11}(1 + l_{11}) + p_{12}(1 + l_{12}) &= (2u_{11} + u_{12} + u_{13})/u_{++} \\
p_{12}(l_{11}l_{12}\lambda_{11} + u_{++}) &= u_{12} \\
(2p_{11}l_{11} + p_{12}l_{12})(u_{++} - l_{11}\lambda_{11}) &= u_{13} \\
p_{12}(1 + l_{11}) + 2p_{22}(1 + l_{12}) &= (u_{12} + 2u_{22} + u_{23})/u_{++} \\
(p_{12}l_{11} + 2p_{22}l_{12})(u_{++} - l_{12}\lambda_{11}) &= u_{23} \\
(2p_{11}l_{11} + p_{12}l_{12})(1 + l_{11}) + (p_{12}l_{11} + 2p_{22}l_{12})(1 + l_{12}) &= (u_{13} + u_{23} + 2u_{33})/u_{++}.
\end{aligned}
$$

This system has six complex solutions for a general data matrix $U$. In the other words, the 6 unknowns $l_{..}, p_{..}$, and $\lambda_{11}$ are algebraic functions of degree 6 in $u_{11}, u_{12}, \ldots, u_{33}$.     □

Here is the symmetric version of Theorem 2, as suggested by Theorem 4:

**Theorem 5** (Draisma and Rodriguez [8])**.** *The ML degree for symmetric $n \times n$-matrices (12) of rank $r$ is equal to the ML degree for symmetric $n \times n$-matrices (12) of rank $n - r + 1$.*

This was stated as a conjecture in the first version of this paper, and proved later in [8].

## 3. Solutions using numerical algebraic geometry

Theorems 1 and 4 document considerable advances relative to the computational results found earlier in [15, §5]. In this project, we used numerical algebraic geometry [5] to compute the ML degrees by solving the local kernel formulation (11) which we explain in this section.

The statistical problem addressed here is to find the global maximum of a likelihood function $\ell_U$ over a matrix model $\mathcal{M}$ given by rank constraints. For this class of problems, the use of numerical algebraic geometry has the following significant advantage over symbolic computations. After having solved the likelihood equations only once, for one generic data matrix $U_0$, all subsequent computations for other data matrices $U$ are much faster. Numerical homotopy continuation will start from the critical points of $\ell_{U_0}$ and transform them into the critical points of $\ell_U$. Intuitively speaking, for a fixed model $\mathcal{M}$, *the homotopy amounts to changing the data.* We believe that our methodology will be useful for a wider range of maximum likelihood problems than those treated here, and we decidedly agree with the statement of Buot and Richards [6, §5] that *"... homotopy continuation algorithms often provide substantial advantages over iterative methods commonly used in statistics".*

We discuss below two options for the preprocessing stage of solving the local kernel formulation (11) for generic $U_0$. The first option is to use a single homotopy built from an upper bound discussed in Section 2, most notably a polyhedral homotopy built from the polyhedral root count. The second option is to use a sequence of homotopies that intersect the hypersurfaces corresponding to each equation, most notably via regeneration [13].

Parallel computation is an essential feature of numerical algebraic geometry. Both preprocessing, by solving a generic data set once, and each subsequent solve for given specific data can be performed in parallel. In our case, we used a 64-bit Linux cluster with 160 processors to perform the computations summarized in Table 2 which tracked each path on a separate processor. For instance, for $(m, n, r) = (4, 5, 2)$, there are 843 paths, to be distributed among the 160 processors. Using adaptive precision [3], this takes 20 seconds while the same computation performed sequentially takes about 20 minutes on a typical laptop.

| $(m, n, r)$ | $(4, 4, 2)$ | $(4, 4, 3)$ | $(4, 5, 2)$ | $(4, 5, 3)$ | $(5, 5, 2)$ | $(5, 5, 4)$ |
|---|---|---|---|---|---|---|
| Preprocessing | 257 | 427 | 1938 | 2902 | 348555 | 146952 |
| Solving | 4 | 4 | 20 | 20 | 83 | 83 |

Table 2: Comparison of running times for preprocessing and subsequent solving (in seconds)

**Example 4.** The following data matrix is attributed to the fictional character DiaNA in [21, Example 1.3]. It represents her alignment of two DNA sequences of length $u_{++} = 40$:

$$U \quad = \quad \begin{pmatrix} 4 & 2 & 2 & 2 \\ 2 & 4 & 2 & 2 \\ 2 & 2 & 4 & 2 \\ 2 & 2 & 2 & 4 \end{pmatrix}.$$

According to Table 2, it took 257 seconds to solve the first instance for $(m, n, r) = (4, 4, 2)$, but now every subsequent run takes only 4 seconds. In that solving step, the integers $u_{ij}$ become parameters over the complex numbers. For DiaNA's data matrix $U$, the 191 complex critical points degenerate to 25 real critical points, each of which is positive, and 166 nonreal critical points. See Theorem 8 for additional information regarding the critical points.     □

Three advantages of the local kernel formulation (11) are that it is a square system with polynomials of degree at most 4, it is sparse in terms of the number of monomials appearing, and it has a natural product structure. These structures are clearly visible from the systems in Example 2, and they are used to derive the smaller upper bounds in Table 1. In what follows, we shall describe our preprocessing and how we can use its output to easily compute all critical points of $\ell_U$ for a given data matrix $U$. We also analyze some specific examples. An introduction to numerical algebraic geometry and homotopy continuation can be found in [22] and more details using `Bertini` to perform these computations appear in the book [5].

For a square polynomial system $F$, *basic homotopy continuation* computes a finite set $\mathcal{S}$ of complex roots of $F$ which contains all isolated roots. Here, "computes $\mathcal{S}$" means numerically computing the coordinates of each point in $\mathcal{S}$, and to be able to approximate these to arbitrary accuracy. Numerical approximations to nonsingular solutions can be certified using the software `alphaCertified` [14]. This certification can also determine if the solution is real or positive. To compute $\mathcal{S}$, we first construct a family of polynomial systems $\mathcal{F}$ containing $F$ and then compute the isolated roots for a sufficiently general $G \in \mathcal{F}$. Finally, one tracks the solution paths starting with the isolated roots as $G$ deforms to $F$ inside $\mathcal{F}$.

Fix $(m, n, r)$ and let $\mathcal{F} := \mathcal{F}_{m,n,r}$ be the family of polynomial systems (11) for $U \in \mathbb{C}^{m \times n}$. The generic root count on $\mathcal{F}$ is the ML degree of $\mathcal{V}_r$. In particular, for any generic $U_0 \in \mathbb{C}^{m \times n}$ the number of roots of the corresponding system $F_{U_0} \in \mathcal{F}$ is the ML degree of $\mathcal{V}_r$. Suppose further that we know the roots of $F_{U_0}$. Then, for any matrix $U \in \mathbb{C}^{m \times n}$, we can compute the isolated roots of the corresponding polynomial system $F_U$ by tracking the ML degree number of solutions paths starting with the roots of $F_{U_0}$ as $U_0$ and $F_{U_0}$ deform to $U$ and $F_U$.

Since the family $\mathcal{F}$ is parameterized by the linear space $\mathbb{C}^{m \times n} \cong \mathbb{R}^{2mn}$, we can connect $U_0$ to $U$ along a line segment. If $U_0$ is not in a sufficiently general position with respect to $U$, e.g., both real, this segment may contain matrices for which the corresponding system has a root count that is different from the ML degree. To avoid this, we apply the *gamma trick* of [20]. For $\gamma \in \mathbb{S}^1 \subset \mathbb{C}^*$, the trick deforms from $U_0$ to $U$ along the arc parameterized by

$$\frac{\gamma t}{1 + (\gamma - 1)t} \cdot U_0 + \frac{1 - t}{1 + (\gamma - 1)t} \cdot U \quad \text{for} \quad t \in [0, 1]. \tag{18}$$

For all but finitely many values $\gamma \in \mathbb{S}^1$, the root count for the corresponding polynomial system along this arc, except possibly at $U$ when $t = 0$, is the ML degree.

We conclude our discussion on deforming from a known set of critical points with a practical issue. Due to choices of affine patches, the local kernel formulation (11), as written, is not

| $(m, n, r)$ | $(3, 3, 2)$ | $(3, 4, 2)$ | $(3, 5, 2)$ | $(4, 4, 2)$ | $(4, 4, 3)$ |
|---|---|---|---|---|---|
| Polyhedral using `PHC` | 4 | 120 | 2017 | 23843 | 1869 |
| Regeneration using `Bertini` | 6 | 61 | 188 | 2348 | 7207 |

Table 3: Running times for preprocessing in serial using double precision (in seconds)

suitable for a nongeneric data matrix $U$. Once given a data matrix $U$, we simply choose random affine patches as in [2]. Let $O_1, O_2 \in \mathbb{R}^{r \times r}$, $O_3 \in \mathbb{R}^{m \times m}$, and $O_4 \in \mathbb{R}^{n \times n}$ be random orthogonal matrices and $L_1$, $P_1$, $R_1$, and $\Lambda$ be as before. Then, we use (11) with

$$
L = O_1 \cdot \begin{pmatrix} L_1 & -I_{m-r} \end{pmatrix} \cdot O_3^T, \quad P = O_3 \cdot \begin{pmatrix} P_1 & P_1 R_1 \\ L_1 P_1 & L_1 P_1 R_1 \end{pmatrix} \cdot O_4^T, \quad \text{and} \quad R = O_4 \cdot \begin{pmatrix} R_1 \\ -I_{n-r} \end{pmatrix} \cdot O_2^T.
$$

The homotopy (18) quickly computes the isolated critical points for any given data matrix $U$ provided that we already know the critical points for a sufficiently general data matrix $U_0$.

We now discuss the two options for *preprocessing* mentioned above, namely polyhedral homotopies and regeneration. A summary of our computations with these two methods, now using serial processing with double precision, are presented in Table 3. The last pair of entries suggest that the two methods exhibit complementary behavior with respect to the duality of Theorem 2. In both cases, 191 roots are found, and these are essentially the same roots, by Theorem 7 below. For instance, using polyhedral homotopy, the rank 2 case can be solved in 1869 seconds and then we may read off the solutions for rank 3 using (19).

The first approach to solve the equations for $U_0$ is to use basic homotopy continuation in the family $\mathcal{P}$ of polynomial systems that arise from some relevant structure. The generic root count on $\mathcal{P}$ constructed from various structures are presented in Table 1. After computing the roots for a general element of $\mathcal{P}$, we return to basic homotopy continuation for computing the roots of $F_{U_0}$. Table 3 summarizes using a polyhedral approach implemented in `PHC` [27] where the family $\mathcal{P}$ is constructed based on the Newton polytopes of the given equations.

The second approach is based on intersecting the given hypersurfaces iteratively. This can be advantageous when the degree of the intersection is significantly less than the product of the degrees. To be explicit, if $\mathcal{S}$ is a pure $k$-dimensional variety ($k > 0$) and $\mathcal{H}$ is a hypersurface, intersection approaches can be advantageous when the degree of the pure $(k-1)$-dimensional part of $\mathcal{S} \cap \mathcal{H}$ is less than $\deg \mathcal{S} \cdot \deg \mathcal{H}$. Regeneration is an intersection approach that builds from a product structure of the given system. We shall now discuss this.

We first consider the classical idea of solving polynomial systems using successive intersections and then discuss how to build from a product structure. Consider $N$ polynomials $f_1, \ldots, f_N$ in $N$ variables, defining hypersurfaces $\mathcal{H}_1, \ldots, \mathcal{H}_N$. One advantage of a square system is that the isolated solutions of $f_1 = \cdots = f_N = 0$ arise by computing the codimension $i$ components of $\mathcal{H}_1 \cap \cdots \cap \mathcal{H}_i$ sequentially for $i = 1, 2, \ldots, N$. In fact, every codimension $i+1$ component of $\mathcal{H}_1 \cap \cdots \cap \mathcal{H}_i \cap \mathcal{H}_{i+1}$ arises as the intersection of a codimension $i$ component $C$ of $\mathcal{H}_1 \cap \cdots \cap \mathcal{H}_i$ and the hypersurface $\mathcal{H}_{i+1}$, where $C$ is not contained in $\mathcal{H}_{i+1}$.

The use of the product structure arises from intersecting an algebraic set of pure codimension $i$ with a linear space of dimension $i$ yielding finitely many points. The first step is a hypersurface intersected with a line. If $\mathcal{L}_2, \ldots, \mathcal{L}_N$ are general hyperplanes, the hypersurface $\mathcal{H}_1$ is represented by the isolated points in $\mathcal{H}_1 \cap \mathcal{L}_2 \cap \cdots \cap \mathcal{L}_N$. Such points can be computed by solving a univariate polynomial, namely $f_1$ restricted to the line $\mathcal{L}_2 \cap \cdots \cap \mathcal{L}_N$. Let $1 \leq i < N$ and $C_i$ be the pure one-dimensional component of $\mathcal{H}_1 \cap \cdots \cap \mathcal{H}_i \cap \mathcal{L}_{i+2} \cap \cdots \cap \mathcal{L}_N$. Now, *basic regeneration* computes $C_i \cap \mathcal{H}_{i+1}$ from $C_i \cap \mathcal{L}_{i+1}$ as follows. Let $\mathcal{M}_1, \ldots, \mathcal{M}_k$ be hyperplanes defined by sufficiently general linear polynomials $\ell_1, \ldots, \ell_k$ that represent a linear product decomposition of $f_{i+1}$. Let $\mathcal{M} = \bigcup_{j=1}^k \mathcal{M}_j$. Basic homotopy continuation computes $C_i \cap \mathcal{M}_j$ from $C_i \cap \mathcal{L}_{i+1}$ for $j = 1, \ldots, k$. Their union is $C_i \cap \mathcal{M}$. Applying basic homotopy continuation once more yields $C_i \cap \mathcal{H}_{i+1}$ by deforming from $C_i \cap \mathcal{M}$.

For the preprocessing approaches above, we can certify that the set of approximations obtained correspond to distinct solutions using `alphaCertified`. At each stage of the regeneration and at the end of the computation, we can perform one additional test to confirm that we have obtained all of the solutions: the trace test [23]. During regeneration, the centroid of the solutions must move linearly as the hyperplane $\mathcal{L}_N$ is moved linearly. Moreover, the centroid of the critical $m \times n$-matrices must move linearly as the data matrix $U$ moves linearly. With these tests, we are able to claim, with high probability, that our initial randomly selected data matrix $U_0$ was sufficiently generic, and Theorems 1 and 4 hold.

After computing the positive critical points for a given data matrix $U$, we identify the local maximizers by analyzing the Hessian of the corresponding Lagrangian function, namely

$$L(P, \lambda) \;=\; \log \ell_U(P) + \sum_{i=1}^{k} \lambda_i g_i(P),$$

where $\mathcal{V}_r$ is defined by the vanishing of the polynomials $g_1, \ldots, g_k$. If $P$ is a critical point of rank $r$, let $\lambda \in \mathbb{C}^k$ be the unique vector such that $\nabla L(P, \lambda) = 0$. Then, $P$ is a local maximizer if the matrix $N^T \cdot HL(P, \lambda) \cdot N$ is negative semidefinite where $HL(P, \lambda)$ is the Hessian of $L$ and the columns of $N$ form a basis for the tangent space of $\mathcal{V}_r \times \mathbb{C}^k$ at $(P, \lambda)$.

In the remainder of this section we present three concrete numerical examples.

**Example 5.** We consider the symmetric matrix model (12) for $n = 3$ with the data

$$u_{11} = 10, \; u_{12} = 9, \; u_{13} = 1, \; u_{22} = 21, \; u_{23} = 3, \; u_{33} = 7.$$

All six critical points of the likelihood function (14) are real and positive. They are

| $p_{11}$ | $p_{12}$ | $p_{13}$ | $p_{22}$ | $p_{23}$ | $p_{33}$ | $\log \ell_U(p)$ |
|---|---|---|---|---|---|---|
| 0.1037 | 0.3623 | 0.0186 | 0.3179 | 0.0607 | 0.1368 | $-82.18102$ |
| 0.1084 | 0.2092 | 0.1623 | 0.3997 | 0.0503 | 0.0702 | $-84.94446$ |
| 0.0945 | 0.2554 | 0.1438 | 0.3781 | 0.4712 | 0.0810 | $-84.99184$ |
| 0.1794 | 0.2152 | 0.0142 | 0.3052 | 0.2333 | 0.0528 | $-85.14678$ |
| 0.1565 | 0.2627 | 0.0125 | 0.2887 | 0.2186 | 0.0609 | $-85.19415$ |
| 0.1636 | 0.1517 | 0.1093 | 0.3629 | 0.1811 | 0.0312 | $-87.95759$ |

The first three points are local maxima in $\Delta_5$ and the last three points are local minima. These six points define an extension of degree 6 over $\mathbb{Q}$. For instance, via `Macaulay 2` [12], the minimal polynomial for the last coordinate is $9528773052286944 p_{33}^6 - 4125267629399052 p_{33}^5 + 713452955656677 p_{33}^4 - 63349419858182 p_{33}^3 + 3049564842009 p_{33}^2 - 75369770028 p_{33} + 744139872$. As we shall see in Proposition 1, the Galois group of this irreducible polynomial is solvable, so we can express each of the coordinates in radicals. The last coordinate, via `RadiRoot` [7], is
$p_{33} = \frac{16427}{227664} + \frac{1}{12}\big(\zeta - \zeta^2\big)\omega_2 - \frac{66004846384302}{19221271018849}\omega_2^2 + \big(\frac{14779904193}{211433981207339}\zeta^2 - \frac{14779904193}{211433981207339}\zeta\big)\omega_1\omega_2^2 + \frac{1}{2}\omega_3$,
where $\zeta$ is a primitive third root of unity, $\omega_1^2 = 94834811/3$, and

$$
\begin{aligned}
\omega_2^3 &= \;\big(\tfrac{5992589425361}{150972770845322208}\zeta - \tfrac{5992589425361}{150972770845322208}\zeta^2\big) + \tfrac{97163}{40083040181952}\omega_1, \\
\omega_3^2 &= \;\tfrac{5006721709}{1248260766912} + \big(\tfrac{212309132509}{4242035935404}\zeta - \tfrac{212309132509}{4242035935404}\zeta^2\big)\omega_2 - \tfrac{2409}{20272573168}\omega_1\omega_2 \\
&\quad - \tfrac{158808750548335}{76885084075396}\omega_2^2 + \big(\tfrac{17063004159}{422867962414678}\zeta^2 - \tfrac{17063004159}{422867962414678}\zeta\big)\omega_1\omega_2^2.
\end{aligned}
$$

We finally note that the six critical points can be matched into three pairs so that (19) holds: the Hadamard product of points 1 and 6 agree with that of points 2 and 5, and that of points 3 and 4. Thus this example illustrates the symmetric matrix version of Theorem 7. $\qquad\square$

**Example 6.** Let $m = 4, n = 5$ and consider the data matrix

$$U \quad = \quad \begin{pmatrix} 2084 & 1 & 1 & 1 & 4 \\ 4 & 23587 & 5 & 3 & 1 \\ 6 & 3 & 41224 & 3 & 2 \\ 4 & 6 & 2 & 8734 & 4 \end{pmatrix}.$$

For $r = 2$ and $r = 3$, this instance has the expected number 843 of distinct complex critical points. In both cases, 555 critical points are real, and 25 of these are positive. Consider the 25 critical points in $\Delta_{19}$. For $r = 2$ precisely seven are local maxima, and for $r = 3$ precisely six are local maxima. We shall list them explicitly in Examples 8 and 9 respectively. □

**Example 7.** Let $m = n = 5$, with the non-symmetric model, and consider the data

$$U \quad = \quad \begin{pmatrix} 2864 & 6 & 6 & 3 & 3 \\ 2 & 7577 & 2 & 2 & 5 \\ 4 & 1 & 7543 & 2 & 4 \\ 5 & 1 & 2 & 3809 & 4 \\ 6 & 2 & 6 & 3 & 5685 \end{pmatrix}.$$

For $r = 2$ and $r = 4$, this instance has the expected number of 6776 distinct complex critical points. In both cases, 1774 of these are real and 90 of these are real and positive. This illustrates the last statement in Theorem 7. The number of local maxima for $r = 2$ equals 15, and the number of local maxima for $r = 4$ equals 6. For $r = 3$, we have 61326 critical points, of which 15450 are real. Of these, 362 are positive and 25 are local maxima. □

## 4. Further results and computations

The numerical algebraic geometry techniques described in Section 3 have the advantage that they permit fast experimentation with non-trivial instances. This led us to a variety of conjectures, including those concerning ML duality. Before we come to our discussion of duality, we briefly state a conjecture regarding the ML degree of $3 \times n$-matrices of rank 2.

**Conjecture 6.** *For $m = 3$ and $n \geq 3$, the ML degree of the variety $\mathcal{V}_2$ equals $2^{n+1} - 6$.*

The first three values already appeared in Theorem 1. We tested this formula by solving the equations of the local kernel formulation (11). This was done independently in Macaulay 2 and Bertini. With these computations, we verified Conjecture 6 up to $n = 10$. This conjecture, if correct, would furnish a simple and natural sequence of models, namely $3 \times n$-matrices of rank 2, whose ML degree grows exponentially in the number of states.

We next formulate a refined version of the duality statement in Theorem 2. Given a data matrix $U$ of format $m \times n$, we write $\Omega_U$ for the $m \times n$-matrix whose $(i, j)$ entry equals

$$\frac{u_{ij} u_{i+} u_{+j}}{(u_{++})^3}.$$

The following statement also appeared as a conjecture in the first version of our paper, and it was proved by Draisma and Rodriguez in their article [8] on maximum likelihood duality.

**Theorem 7** ([8]). *Fix $m \leq n$ and $U$ an $m \times n$-matrix with strictly positive integer entries. There exists a bijection between the complex critical points $P_1, P_2, \ldots, P_s$ of the likelihood function $\ell_U$ on $\mathcal{V}_r$ and the complex critical points $Q_1, Q_2, \ldots, Q_s$ of $\ell_U$ on $\mathcal{V}_{m-r+1}$ such that*

$$P_1 \star Q_1 = P_2 \star Q_2 = \cdots = P_s \star Q_s = \Omega_U. \tag{19}$$

*In particular, this bijection preserves reality, positivity, and rationality of the critical points.*

From the perspective of statistics, this result implies the following striking statement: maximum likelihood estimation for matrices of rank $r$ is exactly the same problem as minimum likelihood estimation for matrices of corank $r - 1$, and vice versa. This refined formulation of the duality statement allows us to improve the speed of MLE by passing to the complementary problem, where it may be easier to solve the likelihood equations. We saw a first instance of this in Section 3 when we discussed the last two columns in Table 3: the two methods give the same set of 191 solutions but the running times are complementary.

**Remark 1.** Equation (19) is trivially satisfied for $r = 1$, where the ML degree is $s = 1$. Here, $P_1$ is the rank one matrix in (22), and $Q_1 = \frac{1}{u_{++}}U$. Clearly, we have $P_1 \star Q_1 = \Omega_U$.    □

We illustrate Theorem 7 for a specific case that has already appeared in the literature [10, 21, 29]. The first assertion in the next theorem resolves [29, Conjecture 11] affirmatively. In their conjecture, Zhu *et al.* [29] had identified the matrix $P(a, b)$ below, and they had asserted that it is the global maximum of the likehood function for the data matrix $U(a, b)$. Note that, for $a = 4$ and $b = 2$, this is the matrix for DiaNA's data in [21, Example 1.16].

**Theorem 8.** *Let $m = n = 4$, $a > b > 0$, and consider the following matrices:*

$$U(a,b) = \begin{bmatrix} a & b & b & b \\ b & a & b & b \\ b & b & a & b \\ b & b & b & a \end{bmatrix} \quad and \quad P(a,b) = \frac{1}{8(a+3b)} \begin{bmatrix} a+b & a+b & 2b & 2b \\ a+b & a+b & 2b & 2b \\ 2b & 2b & a+b & a+b \\ 2b & 2b & a+b & a+b \end{bmatrix}.$$

*The distribution $P(a, b)$ maximizes the likelihood function for the data matrix $U(a, b)$ on $\mathcal{M}_2$.*

*Proof.* This statement is invariant under scaling the vector $(a, b)$. We normalize by taking $4a + 12b = 16$. Then $b = (4 - a)/3$ and $a$ ranges in the open interval defined by $1 < a < 4$. For each such $a$, the likelihood function $\ell_{U(a,b)}$ has exactly 25 positive critical points in the rank 2 model $\mathcal{M}_2$, with the maximum value occurring at $P(a, b)$. This statement was shown using the following method and its illustration in Figure 1.
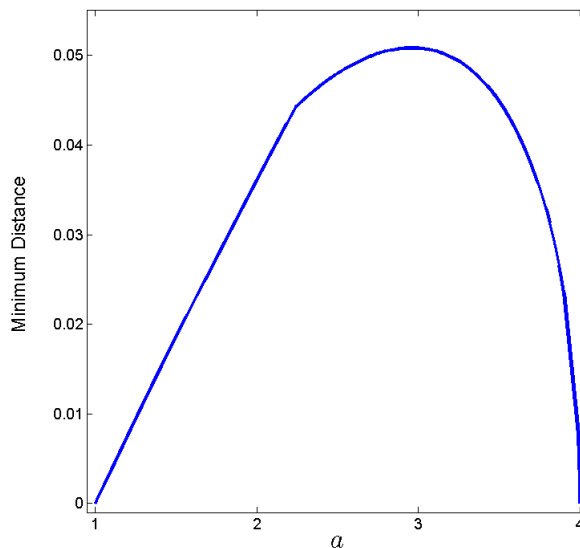


Figure 1: Minimum pairwise distance and lower bound (20) as a function of $a$

First, we selected $a = 2$ and computed the 191 critical points using `Bertini`. From these, `alphaCertified` proved that exactly 25 are real and, using the computed error bounds, it verified that all lie in $\Delta_{15}$. We then expressed these real solutions as rational functions in $a$ and $b$ to show that all 25 real solutions remain positive for all $a > b > 0$. The critical points fall

into four symmetry classes of size 6, 12, 4, and 3. Representatives of these classes are

$$
X_1 = \frac{1}{16} \begin{bmatrix} 1 & 1 & 1 & 1 \\ 1 & 1 & 1 & 1 \\ 1 & 1 & \frac{2a}{a+b} & \frac{2b}{a+b} \\ 1 & 1 & \frac{2b}{a+b} & \frac{2a}{a+b} \end{bmatrix}, \qquad
X_2 = \frac{1}{32(a+2b)} \begin{bmatrix} 2a+4b & 2a+4b & 2a+4b & 2a+4b \\ 2a+4b & 6a & 6b & 6b \\ 2a+4b & 6b & 3a+3b & 3a+3b \\ 2a+4b & 6b & 3a+3b & 3a+3b \end{bmatrix},
$$

$$
X_3 = \frac{1}{12(a+3b)} \begin{bmatrix} 3a & 3b & 3b & 3b \\ 3b & a+2b & a+2b & a+2b \\ 3b & a+2b & a+2b & a+2b \\ 3b & a+2b & a+2b & a+2b \end{bmatrix}, \qquad \text{and} \quad X_4 = P(a,b).
$$

Using calculus, one can prove that $\log \ell_U(X_i) < \log \ell_U(X_{i+1})$ for $i = 1, 2, 3$.

All that remains is to show that the 191 solutions remain distinct on $1 < a < 4$ (with some coalesce at the boundary). The function mapping $a$ to the minimum of the pairwise distances between the critical points is a piecewise smooth function. It is depicted in Figure 1. By tracking the homotopy paths as $a$ changes from 2 to 1 and from 2 to 4, we are able to determine that this function is nowhere zero on the open interval $(1, 4)$. Additionally, by analyzing the solutions using [1], a lower bound on this minimum pairwise distance function is

$$
\min \left\{ \frac{(a-1)\sqrt{a^2+17}}{12(a+8)}, \quad \frac{a+2-\sqrt{(a-1)(a-4)}}{48} - \frac{3(a^2-12a-16)+\sqrt{6(a-1)(a-4)(a^2-16a+96)}}{16(a+8)(a-10)} \right\} \tag{20}
$$

which is also depicted in Figure 1. The first term of this minimum arises from $X_2$ and a member of the $X_3$ family which is equal to the minimum pairwise distances for values of $a$ near 1. The second term arises from comparing the $(1, 1)$ entries of critical points. In short, all of the solutions remain distinct on $1 < a < 4$ and this establishes [29, Conjecture 11].

We checked the duality statement in Theorem 7 by performing the same computation for $m = n = 4$ and $r = 3$. We followed the 191 paths in the deformation from a general $U_0$ to a general $U(a, b)$. Using Bertini, we found that 12 endpoints had rank 2 while the other 179 had the expected rank of 3. Moving the other 179 solutions to $a = 2$ produced 179 distinct complex solutions that remain distinct and retain rank 3 on $(1, 4)$. Using the same certification process as above, precisely 25 are positive. These critical points of $\mathcal{M}_3$ form four symmetry classes having the same sizes $6, 12, 4$, and 3 as above, with representatives:

$$
Y_1 = \frac{1}{8(a+3b)} \begin{bmatrix} 2a & 2b & 2b & 2b \\ 2b & 2a & 2b & 2b \\ 2b & 2b & a+b & a+b \\ 2b & 2b & a+b & a+b \end{bmatrix}, \qquad
Y_2 = \frac{1}{12(a+3b)} \begin{bmatrix} 3a & 3b & 3b & 3b \\ 3b & a+2b & a+2b & a+2b \\ 3b & a+2b & \frac{2a(a+2b)}{a+b} & \frac{2b(a+2b)}{a+b} \\ 3b & a+2b & \frac{2b(a+2b)}{a+b} & \frac{2a(a+2b)}{a+b} \end{bmatrix},
$$

$$
Y_3 = \frac{1}{16(a+2b)} \begin{bmatrix} a+2b & a+2b & a+2b & a+2b \\ a+2b & 3a & 3b & 3b \\ a+2b & 3b & 3a & 3b \\ a+2b & 3b & 3b & 3a \end{bmatrix}, \quad
Y_4 = \frac{1}{16(a+b)} \begin{bmatrix} 2a & 2b & a+b & a+b \\ 2b & 2a & a+b & a+b \\ a+b & a+b & 2a & 2b \\ a+b & a+b & 2b & 2a \end{bmatrix}.
$$

The matrices are now sorted by decreasing value of $\ell_{U(a,b)}$, so the first matrix $Y_1$ is the MLE. Our real positive critical points satisfy the desired duality relation. Namely, we have

$$
X_1 \star Y_1 = X_2 \star Y_2 = X_3 \star Y_3 = X_4 \star Y_4 = \frac{1}{64(a+3b)} U(a, b) =: \Omega_U.
$$

We verified the same for the complex solutions.

When Theorem 7 was still a conjecture, we verified it for randomly selected data matrices with i.i.d. entries sampled from the uniform distribution on $[0, 1]$. After generating a random matrix, we verified equation (19) using the critical points computed by homotopy continuation. For $m = n = 3$ and $r = 2$, we verified (19) for 50000 instances. Additionally, for $m = n = 4$

and $r = 2$, we verified (19) for 10000 instances. We also did this for a handful of $4 \times 5$ instances (such as Example 6) and $5 \times 5$ instances (such as Example 7). The user can find `Macaulay 2` code, which uses the emerging `Bertini.m2` package, to perform more numerical experiments at `http://www.math.ncsu.edu/~jdhauens/MLE`.

Theorem 7 and its analogue for symmetric matrices is particularly interesting in the special case when $m = n = 2r - 1$. Here we have an involution on the set of critical points of $\ell_U$ on $\mathcal{V}_r$ which has the following property. If $P_1, P_2, \ldots, P_s$ are the positive critical points in the model $\mathcal{M}_r$, ordered by increasing value of the log-likelihood function, then

$$\ell_U(P_1) + \ell_U(P_s) = \ell_U(P_2) + \ell_U(P_{s-1}) = \cdots = \ell_U(P_{\lceil s/2 \rceil}) + \ell_U(P_{\lfloor s/2 \rfloor}).$$

The identity (19) implies that Galois group which permutes the set of critical points is considerably smaller than the full symmetric group on these points. We shall demonstrate this for $n = 3$. What follows will explain the solutions in radicals seen in Example 5.

Let $\mathbb{Q}(U)$ denote the field of rational functions in entries of an indeterminate data matrix $U$, and let $K$ denote the algebraic extension of $\mathbb{Q}(U)$ that is defined by adjoining all solutions of the likelihood equations. Thus the degree of the extension $K/\mathbb{Q}(U)$ is the ML degree. We are interested in the Galois group $G = \mathrm{Gal}(K, \mathbb{Q}(U))$ of this algebraic extension. This Galois group is a subgroup of the full symmetric group $S_M$ where $M$ is the ML degree.

The following result was found by explicit computations using `maple` and `Sage` [25].

**Proposition 1.** *The Galois group for MLE on $3 \times 3$-matrices (1) of rank 2 is a subgroup of order* 1920 *in $S_{10}$. As an abstract group, it is the semidirect product of $S_5$ and $(\mathbb{Z}_2)^4$. The Galois group for MLE on symmetric $3 \times 3$-matrices (12) of rank 2 is a subgroup of order 24 in $S_5$. As an abstract group, it is the symmetric group $S_4$. So, in the latter case, the six critical points of the likelihood function can be written in radicals in $u_{11}, u_{12}, u_{13}, u_{22}, u_{23}, u_{33}$.*

We close this section with an important observation that is implied by the various polynomial formulations of our problem, but which had not been explicitly stated in Section 2.

**Remark 2.** *Every complex critical point $P$ of the likelihood function $\ell_U$ on $\mathcal{V}_r$ satisfies*

$$p_{i+} = \frac{u_{i+}}{u_{++}} \;\; for \; i = 1, \ldots, m \qquad and \qquad p_{+j} = \frac{u_{+j}}{u_{++}} \;\; for \; j = 1, \ldots, n.$$

The analogous identities hold for any statistical model that is *toric* in the sense of [21]. Namely, the critical points of the likelihood function on any secant variety of a toric variety have the sufficient statistics of the given data in the toric model. This fact seems relevant for the topological underpinnings of ML duality. One is tempted to speculate that some version of Theorems 2, 5, and 7 might be true for other classes of toric models.

## 5. Rank versus non-negative rank

In the previous sections, we developed accurate methods for finding the global maximum of a likelihood function $\ell_U$ over non-negative matrices $P$ of rank $r$ whose entries sum to 1. Unfortunately, this is not quite the problem most practitioners and users of statistics would actually be interested in. Rather than restricting the rank of a probability table (1), it is the *non-negative rank* that is more relevant for applications. In this section we discuss this.

Let $\mathrm{Mix}_r$ denote the subset of $\Delta_{mn-1}$ that comprises all the mixtures of $r$ independent distributions. In statistics, this is the archetype of a latent variable model, or hidden variable model. Mathematically, we can define the *mixture model* $\mathrm{Mix}_r$ as the set of all matrices

$$P = A \cdot \Lambda \cdot B, \tag{21}$$

where $A$ is a non-negative $m \times r$-matrix whose columns sum to 1, $\Lambda$ is an $r \times r$ diagonal matrix whose diagonal entries are non-negative and sum to 1, and $B$ is a non-negative $r \times n$-matrix

whose rows sum to 1. The *rank-constrained model* $\mathcal{M}_r = \mathcal{V}_r \cap \Delta_{mn-1}$ we discussed above is an algebraic relaxation of the mixture model $\mathrm{Mix}_r$. This can be made precise as follows:

**Proposition 2.** *The rank-constrained model $\mathcal{M}_r$ is the Zariski closure of the mixture model $\mathrm{Mix}_r$ inside the simplex $\Delta_{mn-1}$. If $r \leq 2$ then $\mathrm{Mix}_r = \mathcal{M}_r$. If $r \geq 3$ then $\mathrm{Mix}_r \subsetneq \mathcal{M}_r$.*

*Proof.* See Example 4.1.2, Example 4.1.4 and Proposition 4.1.6 in [9]. That book refers to secant varieties of Segre varieties, tensors of any format, and joint distributions of any number of random variables. Here we only need the case of matrices and two random variables.

Our model $\mathcal{M}_r$ is the set of all distributions $P$ of rank at most $r$, while $\mathrm{Mix}_r$ is the set of all distributions $P$ of non-negative rank at most $r$. Having non-negative rank $\leq r$ means that $P = A' \cdot B'$ for some non-negative matrices where $A'$ has $r$ columns and $B'$ has $r$ rows. Any such factorization can be transformed into the particular form (21) which identifies the statistical parameters. For further information on these two models see [10, 19, 21].

Understanding the inclusion of $\mathrm{Mix}_r$ inside $\mathcal{M}_r$ becomes crucial when comparing different methodologies for maximum likelihood estimation. We used `Bertini` to compute all critical points of the likelihood function $\ell_U$ on $\mathcal{M}_r$, with the aim of identifying the global maximum $\widehat{P}$ of $\ell_U$ over $\mathcal{M}_r$. This assumes that $\widehat{P}$ is strictly positive. This is usually the case when $U$ is strictly positive. The standard method used by statisticians is to run the *EM algorithm* in the space of model parameters $(A, \Lambda, B)$. This results in a local maximum $(\widehat{A}, \widehat{\Lambda}, \widehat{B})$ of the likelihood function expressed in terms of the parameters. The fact that $\mathcal{M}_r$ is the Zariski closure of the mixture model $\mathrm{Mix}_r$ in the simplex $\Delta_{mn-1}$ has the following consequence:

**Corollary 1.** *Let $\widehat{P}_1, \ldots, \widehat{P}_s$ be the local maxima in $\mathcal{M}_r$ of the likelihood function $\ell_U$. If a matrix $\widehat{P}_i$ has non-negative rank at most $r$ then $\widehat{P}_i$ lies in $\mathrm{Mix}_r$ and matching parameters $(\widehat{A}_i, \widehat{\Lambda}_i, \widehat{B}_i)$ can found by solving (21). If all matrices $\widehat{P}_i$ have non-negative rank strictly larger than $r$ then $\ell_U$ attains its maximum over $\mathrm{Mix}_r$ on the topological boundary $\partial \mathrm{Mix}_r$.*

*Proof.* The second sentence holds because every matrix $P \in \Delta_{mn-1}$ of non-negative rank $\leq r$ admits a factorization of the special form (21). Indeed, if $P = A' \cdot B'$ is any non-negative factorization then we first scale the rows of $A'$ to get a matrix $A$ with row sums equal to 1, and we adjust the second matrix so that $P = A \cdot B''$. Now let $\Lambda$ be the diagonal matrix whose entries are the column sums of $B''$ and set $B = \Lambda^{-1} B''$. Then $P = A \cdot \Lambda \cdot B$.

For the third sentence, suppose $\ell_U$ has its maximum over $\mathrm{Mix}_r$ at a point $\widehat{P}$ in $\mathrm{Mix}_r \backslash \partial \mathrm{Mix}_r$. Then $\widehat{P}$ is also a local maximum of $\ell_U$ on $\mathcal{M}_r$. Thus $\widehat{P}$ will be found by solving the critical equations for $\ell_U$ on $\mathcal{V}_r$. The matrix $\widehat{P}$ is an element of $\{\widehat{P}_1, \ldots, \widehat{P}_s\}$. Hence, this set contains a matrix of non-negative rank $\leq r$. This proves the contrapositive of the assertion.

We shall now discuss the exact solution of the MLE problem for the mixture model $\mathrm{Mix}_r$. Let us start with the low rank cases. The given input is a data matrix $U$ as in (2).

If $r = 1$ then the likelihood function $\ell_U$ has a unique critical point. Let $u_{*+}$ be the column vector of row sums of $U$, and let $u_{+*}$ be the row vector of column sums of $U$. Then

$$\widehat{P} \;\;=\;\; \frac{1}{(u_{++})^2} \cdot u_{*+} \cdot u_{+*}. \tag{22}$$

If $r \geq 2$ then we compute the set $\{\widehat{P}_1, \ldots, \widehat{P}_s\}$ of all local maxima of the likelihood function $\ell_U$ on the model $\mathcal{M}_r$. This is done using the numerical algebraic geometry methods described in Section 3, by solving the likelihood equations (11) for the determinantal variety $\mathcal{V}_r$.

If $r = 2$ then every matrix $\widehat{P}_i$ has non-negative rank $\leq 2$. We therefore select the matrix whose likelihood value $\ell_U(\widehat{P}_i)$ is maximal. Then $\widehat{P}_i$ solves the MLE problem for $\mathrm{Mix}_2 = \mathcal{M}_2$.

**Example 8.** We experimented with the EM Algorithm for $r = 2$, as in [21, §1.3], on the $4 \times 5$ data matrix $U$ discussed in Example 6. We ran 10000 iterations with starting points $(A, \Lambda, B)$ sampled from the uniform distribution on the 15-dimensional parameter polytope

$$(\Delta_3 \times \Delta_3) \times \Delta_1 \times (\Delta_4 \times \Delta_4).$$

From these 10000 runs of the EM algorithm we obtained the following seven local maxima:

2643 occurrences:
$$\begin{bmatrix} 0.001678 & 0.01892 & 0.00001325 & 0.007008 & 0.00000722 \\ 0.01894 & 0.2136 & 0.00006605 & 0.07912 & 0.00008149 \\ 0.00007930 & 0.00003964 & 0.5447 & 0.00003964 & 0.00002643 \\ 0.007023 & 0.07921 & 0.00002643 & 0.02933 & 0.00003021 \end{bmatrix} \quad \log(\ell_U) = -105973.49$$

2044 occurrences:
$$\begin{bmatrix} 0.001332 & 0.00001777 & 0.02627 & 0.00000792 & 0.00000382 \\ 0.00007696 & 0.2274 & 0.00006503 & 0.08423 & 0.00004823 \\ 0.02628 & 0.00003913 & 0.5185 & 0.00004103 & 0.00007542 \\ 0.00002871 & 0.08432 & 0.00002762 & 0.03123 & 0.00001788 \end{bmatrix} \quad \log(\ell_U) = -106487.35$$

1897 occurrences:
$$\begin{bmatrix} 0.002245 & 0.02536 & 0.00001725 & 0.000006332 & 0.000005379 \\ 0.02535 & 0.2863 & 0.00006471 & 0.00004393 & 0.00006072 \\ 0.00009818 & 0.00003897 & 0.4495 & 0.09525 & 0.00006537 \\ 0.00002773 & 0.00008630 & 0.09530 & 0.02020 & 0.00001388 \end{bmatrix} \quad \log(\ell_U) = -109697.04$$

1688 occurrences:
$$\begin{bmatrix} 0.001111 & 0.00001327 & 0.02187 & 0.004634 & 0.000005304 \\ 0.00005289 & 0.3117 & 0.00006605 & 0.00003968 & 0.00001322 \\ 0.02191 & 0.00003963 & 0.4314 & 0.09144 & 0.0001046 \\ 0.004647 & 0.00007931 & 0.09148 & 0.01939 & 0.00002219 \end{bmatrix} \quad \log(\ell_U) = -111172.67$$

1106 occurrences:
$$\begin{bmatrix} 0.005321 & 0.00002006 & 0.00001106 & 0.02226 & 0.00002038 \\ 0.00005070 & 0.1135 & 0.1983 & 0.00004009 & 0.00001444 \\ 0.00008126 & 0.1983 & 0.3465 & 0.00003939 & 0.00002520 \\ 0.02227 & 0.00007333 & 0.00002771 & 0.09316 & 0.00008532 \end{bmatrix} \quad \log(\ell_U) = -127069.50$$

529 occurrences:
$$\begin{bmatrix} 0.0008641 & 0.009735 & 0.01701 & 0.00001350 & 0.00000289 \\ 0.009756 & 0.1099 & 0.1921 & 0.00003965 & 0.00003259 \\ 0.01705 & 0.1921 & 0.3357 & 0.00003959 & 0.00005693 \\ 0.00005301 & 0.00007930 & 0.00002642 & 0.1154 & 0.00005294 \end{bmatrix} \quad \log(\ell_U) = -131013.73$$

93 occurrences:
$$\begin{bmatrix} 0.02754 & 0.00001320 & 0.00001319 & 0.00001334 & 0.00005311 \\ 0.00005280 & 0.09999 & 0.1747 & 0.03704 & 0.00002957 \\ 0.00007916 & 0.1747 & 0.3053 & 0.06472 & 0.00005164 \\ 0.00005339 & 0.03706 & 0.06476 & 0.01373 & 0.00001102 \end{bmatrix} \quad \log(\ell_U) = -148501.63$$

The first matrix is the global maximum, and it was the output in 2643 of our 10000 runs. Note that the ordering by objective function value agrees with the ordering by occurrence. We know from Example 6 that $\Delta_{19}$ contains 7 local maxima, and hence our EM experiment found them all. Each of the 7 matrices above has both rank and non-negative rank $r = 2$.                $\square$

If $r \geq 3$ then the situation is more challenging. To begin with, we need a method for testing whether a matrix has non-negative rank $\leq r$. Recent work by Moitra [18] shows that the computational complexity of this problem is lower than one might fear at first glance.

So, let us assume for now that this problem has been solved and we have an algorithm to decide quickly whether any of the matrices $\widehat{P}_i$ has non-negative rank $r$. If so, we pick among them the matrix $\widehat{P}_i$ of largest $\ell_U$-value. This matrix is now a candidate for the MLE on $\mathrm{Mix}_r$. But it may not actually be the MLE because the global maximum of the likelihood function $\ell_U$ may be attained on the boundary $\partial \mathrm{Mix}_r$. Furthermore, it is quite possible that none of the critical points in $\{\widehat{P}_1, \ldots, \widehat{P}_s\}$ lies in $\mathrm{Mix}_r$. Then, according to the third sentence of Corollary 1, the MLE in the mixture model $\mathrm{Mix}_r$ necessarily lies in the boundary $\partial \mathrm{Mix}_r$.

Our discussion implies that, in order to perform exact maximum likelihood estimation for the mixture model, we need to have an exact algebraic description of $\partial \mathrm{Mix}_r$. Specifically, we must determine the polynomial equations that cut out the various irreducible components of the Zariski closure of $\partial \mathrm{Mix}_r$ as a subvariety of $\mathbb{P}^{mn-1}$. For each of these components, and the various strata where they intersect, we then need to compute the ML degree. That list of further ML degrees, combined with the value for $\mathcal{V}_r$ in Theorem 1, describes the true intrinsic algebraic complexity of the MLE $\widehat{P}$ as a piecewise algebraic function of the data $U$.

To be even more ambitious, we could ask for an exact semi-algebraic description of the set $\mathrm{Mix}_r$. Namely, what we seek is a Boolean combination of polynomial inequalities in the unknowns $p_{ij}$ that characterize $\mathrm{Mix}_r$ as a subset of $\mathcal{V}_r \cap \Delta_{mn-1}$. Finding such a description is an open problem, even in the small cases that are covered by Theorem 1. We believe that it might be possible to resolve the problem for these cases, where $(m, n, r)$ ranges from $(4, 4, 3)$ to $(5, 5, 4)$, using the techniques developed by Mond, Smith, and van Straten in [19].

We illustrate the proposed approach for the first interesting case $(m, n, r) = (4, 4, 3)$. Components of $\partial \mathrm{Mix}_3$ correspond to different labelings of the configurations in [19, Figure 9]. Using the translations (seen in [19, §2]) between non-negative factorizations (21) and nested polygons, one of the labelings of [19, Figure 9 (a)] corresponds to the factorization

$$
\begin{pmatrix}
p_{11} & p_{12} & p_{13} & p_{14} \\
p_{21} & p_{22} & p_{23} & p_{24} \\
p_{31} & p_{32} & p_{33} & p_{34} \\
p_{41} & p_{42} & p_{43} & p_{44}
\end{pmatrix}
=
\begin{pmatrix}
0 & a_{12} & a_{13} \\
0 & a_{22} & a_{23} \\
a_{31} & 0 & a_{33} \\
a_{41} & a_{42} & 0
\end{pmatrix}
\cdot
\begin{pmatrix}
0 & b_{12} & b_{13} & b_{14} \\
b_{21} & 0 & b_{23} & b_{24} \\
b_{31} & b_{32} & b_{33} & 0
\end{pmatrix}.
\tag{23}
$$

This equation parametrizes an irreducible divisor in the 14-dimensional variety $\mathcal{V}_3 \subset \mathbb{P}^{15}$. That divisor is one of the irreducible components of the algebraic boundary of $\mathcal{M}_3$. The corresponding prime ideal of height 2 in $\mathbb{Q}[p_{11}, \ldots, p_{44}]$ is obtained by eliminating the 17 unknowns $a_{ij}$ and $b_{ij}$ from the 16 scalar equations in (23). We find that this ideal is generated by the $4 \times 4$-determinant that defines $\mathcal{V}_3$ together with four sextics such as

$$
\begin{aligned}
& p_{11}p_{21}p_{22}p_{32}p_{33}p_{43} - p_{11}p_{21}p_{22}p_{33}^2p_{42} - p_{11}p_{21}p_{23}p_{32}^2p_{43} + p_{11}p_{21}p_{23}p_{32}p_{33}p_{42} - p_{11}p_{22}^2p_{31}p_{33}p_{43} \\
& + p_{11}p_{22}p_{23}p_{31}p_{32}p_{43} + p_{11}p_{22}p_{23}p_{31}p_{33}p_{42} - p_{11}p_{23}^2p_{31}p_{32}p_{42} + p_{12}p_{21}p_{22}p_{33}^2p_{41} - p_{12}p_{21}p_{23}p_{32}p_{33}p_{41} \\
& - p_{12}p_{22}p_{23}p_{31}p_{33}p_{41} + p_{12}p_{23}^2p_{31}p_{32}p_{41} + p_{13}p_{21}^2p_{32}^2p_{43} - p_{13}p_{21}^2p_{32}p_{33}p_{42} - 2p_{13}p_{21}p_{22}p_{31}p_{32}p_{43} \\
& + p_{13}p_{21}p_{22}p_{31}p_{33}p_{42} + p_{13}p_{21}p_{23}p_{31}p_{32}p_{42} + p_{13}p_{22}^2p_{31}^2p_{43} - p_{13}p_{22}p_{23}p_{31}^2p_{42}.
\end{aligned}
$$

What needs to be studied now is the ML degree of this codimension 2 subvariety of $\mathbb{P}^{15}$, and the approach of [16] would lead us to look at the topology of the associated very affine variety.

Described above is the geometry of the MLE problem for the mixture model $\mathrm{Mix}_r$ regarded as a subset of the ambient simplex $\Delta_{mn-1}$. Statisticians, on the other hand, are more accustomed to working in the space of model parameters, which is the product of simplices

$$
(\Delta_{m-1})^r \times \Delta_{r-1} \times (\Delta_{n-1})^r.
\tag{24}
$$

Here our parameters are $(A, \Lambda, B)$. The model $\mathrm{Mix}_r$ is the image of this parameter space in $\Delta_{mn-1}$ under the map (21). That parametrization is very far from identifiable. The reason is that the fibers of $(A, \Lambda, B) \mapsto P$ are semi-algebraic sets of possibly large dimension. In fact, the whole point of the paper [19] is to study the topology of these fibers as $P$ varies.

The expectation-maximization (EM) algorithm is the local method of choice for finding the MLE on the mixture model $\mathrm{Mix}_r$. Our readers might enjoy the exposition given in [21, §1.3]. We emphasize that the EM algorithm operates entirely in the parameter space (24). The likelihood function $\ell_U$ pulls back to a function on the interior of (24). The EM algorithm is an iterative method that converges to a critical point of that function, and, under some mild regularity hypotheses, that critical point $(\widehat{A}, \widehat{\Lambda}, \widehat{B})$ is then a local maximum. The image $\widehat{P}$ of the point in $\mathrm{Mix}_r$ is then a candidate for the global maximum of $\ell_U$ on $\mathrm{Mix}_r$.

**Example 9.** We tried the EM Algorithm also for $r = 3$ on the $4 \times 5$ data matrix $U$ in Examples 6 and 8. We ran 10000 iterations with starting points sampled from the uniform distribution on the 23-dimensional parameter polytope $(\Delta_3)^3 \times \Delta_2 \times (\Delta_4)^3$. From these 10000 runs of the EM algorithm, 9997 converged to one of eight local maxima. Three of the runs led to other fixed points. The following six local maxima are precisely the solutions already found in Example 6. We note that, in this particular instance, it happened that all local maxima in the rank model

$\mathcal{M}_3$ actually lie in Mix$_3$, i.e. they have non-negative rank 3:

3521 occurrences: $\begin{bmatrix} 0.005321 & 0.00001322 & 0.00001322 & 0.02226 & 0.00002039 \\ 0.00005285 & 0.3117 & 0.00006607 & 0.0003964 & 0.00001321 \\ 0.00007929 & 0.00003964 & 0.5447 & 0.00003964 & 0.00002643 \\ 0.02227 & 0.00007927 & 0.00002642 & 0.09316 & 0.00008532 \end{bmatrix}$ $\log(\ell_U) = -84649.67679$

2293 occurrences: $\begin{bmatrix} 0.002244 & 0.02535 & 0.00001324 & 0.00001333 & 0.0000054 \\ 0.02535 & 0.2863 & 0.00006606 & 0.00003961 & 0.00006065 \\ 0.00007929 & 0.00003963 & 0.5447 & 0.00003964 & 0.00002643 \\ 0.00005291 & 0.00007928 & 0.00002643 & 0.1154 & 0.00005289 \end{bmatrix}$ $\log(\ell_U) = -86583.69000$

1678 occurrences: $\begin{bmatrix} 0.001332 & 0.00001326 & 0.02627 & 0.00001341 & 0.0000038 \\ 0.00005289 & 0.3117 & 0.00006607 & 0.00003964 & 0.00001322 \\ 0.02628 & 0.00003963 & 0.5185 & 0.00003961 & 0.00007538 \\ 0.00005296 & 0.00007928 & 0.00002642 & 0.1154 & 0.00005292 \end{bmatrix}$ $\log(\ell_U) = -87698.20128$

1320 occurrences: $\begin{bmatrix} 0.02754 & 0.00001320 & 0.00001321 & 0.00001326 & 0.00005298 \\ 0.00005277 & 0.2274 & 0.00006606 & 0.08423 & 0.00004806 \\ 0.00007928 & 0.00003964 & 0.5447 & 0.00003964 & 0.00002643 \\ 0.00005310 & 0.08430 & 0.00002643 & 0.03122 & 0.00001788 \end{bmatrix}$ $\log(\ell_U) = -98171.25551$

576 occurrences: $\begin{bmatrix} 0.02754 & 0.00001321 & 0.00001320 & 0.00001330 & 0.00005305 \\ 0.00005285 & 0.3117 & 0.00006605 & 0.00003968 & 0.00001322 \\ 0.00007916 & 0.00003964 & 0.4495 & 0.09526 & 0.00006519 \\ 0.00005324 & 0.00007932 & 0.09528 & 0.02019 & 0.00001389 \end{bmatrix}$ $\log(\ell_U) = -102495.4349$

68 occurrences: $\begin{bmatrix} 0.02754 & 0.00001322 & 0.00001321 & 0.00001321 & 0.00005285 \\ 0.00005287 & 0.1135 & 0.1983 & 0.00003968 & 0.00001444 \\ 0.00007927 & 0.1983 & 0.3465 & 0.00003962 & 0.00002520 \\ 0.00005285 & 0.00007930 & 0.00002642 & 0.1154 & 0.00005285 \end{bmatrix}$ $\log(\ell_U) = -121802.8945$

In addition, our runs of the EM algorithm discovered the two local maxima

488 occurrences: $\begin{bmatrix} 0.001678 & 0.01892 & 0.00001325 & 0.007008 & 0.0000072 \\ 0.01894 & 0.2136 & 0.00006605 & 0.07912 & 0.00008149 \\ 0.00007930 & 0.00003964 & 0.5447 & 0.00003964 & 0.00002643 \\ 0.007023 & 0.07921 & 0.00002643 & 0.02933 & 0.00003021 \end{bmatrix}$ $\log(\ell_U) = -105973.4859$

53 occurrences: $\begin{bmatrix} 0.001111 & 0.00001341 & 0.02187 & 0.004634 & 0.0000053 \\ 0.00005299 & 0.3117 & 0.00006602 & 0.00003976 & 0.00001324 \\ 0.02191 & 0.00003960 & 0.4314 & 0.09144 & 0.0001046 \\ 0.004647 & 0.00007935 & 0.09148 & 0.01939 & 0.00002219 \end{bmatrix}$ $\log(\ell_U) = -111172.6663$

These do not satisfy the likelihood equations. They are located on the boundary of Mix$_3$.   $\square$

A forthcoming paper by Kaie Kubjas, Elina Robeva and Bernd Sturmfels [17] will analyze the (algebraic) geometry of the EM algorithm, with focus on the small cases of Theorem 1. Comparison with the methods introduced in this paper opens up the possibility of characterizing conditions under which EM finds the global maximum, as it did in Example 9.

## Acknowledgements

## References

[1] D.J. Bates, J.D. Hauenstein, T.M. McCoy, C. Peterson, and A.J. Sommese: *Recovering exact results from inexact numerical data in algebraic geometry*, Experimental Mathematics (2013), to appear.

[2] D.J. Bates, J.D. Hauenstein, C. Peterson, and A.J. Sommese: *Numerical decomposition of the rank-deficiency set of a matrix of multivariate polynomials*, in "Approximate Commutative Algebra" (eds. L. Robbiano and J. Abbott), Texts and Monographs in Symbolic Computation, Springer, Vienna, 2010, pp. 55–77.

[3] D.J. Bates, J.D. Hauenstein, A.J. Sommese, and C.W. Wampler: *Adaptive multiprecision path tracking*, SIAM J. Numer. Anal. **46** (2008) 722–746.

[4] D.J. Bates, J.D. Hauenstein, A.J. Sommese, and C.W. Wampler: *Bertini: Software for Numerical Algebraic Geometry*, `bertini.nd.edu`, 2006.

[5] D.J. Bates, J.D. Hauenstein, A.J. Sommese, and C.W. Wampler: *Numerically Solving Polynomial Systems with the Software Package Bertini*, SIAM, 2013.

[6] M. Buot and D. Richards: *Counting and locating the solutions of polynomial systems of maximum likelihood equations*, J. Symbolic Computation **41** (2006) 234–244.

[7] A. Distler: *RadiRoot: roots of a polynomial as radicals – a GAP package*, version 2.6, `www.icm.tu-bs.de/ag_algebra/software/radiroot`, 2011.

[8] J. Draisma and J. Rodriguez: *Maximum likelihood duality for determinantal varieties*, International Mathematics Research Notices **18** (2013).

[9] M. Drton, B. Sturmfels and S. Sullivant: *Lectures on Algebraic Statistics*, Oberwolfach Seminars, Vol 39, Birkhäuser, Basel, 2009.

[10] S. Fienberg, P. Hersh, A. Rinaldo and Z. Yi: *Maximum likelihood estimation in latent class models for contingency table data*, Algebraic and Geometric Methods in Statistics, 27–62, Cambridge University Press, 2010.

[11] T. Gao, T.Y. Li, and M. Wu: *Algorithm 846: MixedVol: a software package for mixed-volume computation*, ACM Trans. Math. Software **31** (2005) 555–560.

[12] D.R Grayson and M.E. Stillman: *Macaulay2, a software system for research in algebraic geometry*, `www.math.uiuc.edu/Macaulay2`.

[13] J.D. Hauenstein, A.J. Sommese, and C.W. Wampler: *Regeneration homotopies for solving systems of polynomials*, Math. Comp. **80** (2011) 345–377.

[14] J.D. Hauenstein and F. Sottile: *Algorithm 921: alphaCertified: Certifying solutions to polynomial systems*, ACM Trans. Math. Software **38** (2012) 28.

[15] S. Hoşten, A. Khetan and B. Sturmfels: *Solving the likelihood equations*, Foundations of Computational Mathematics **5** (2005) 389–407.

[16] J. Huh: *The maximum likelihood degree of a very affine variety*, Compositio Mathematica **149** (2013) 1245–1266.

[17] K. Kubjas, E. Robeva and B. Sturmfels: Nonnegative matrix rank and the EM algorithm, in preparation.

[18] A. Moitra: *A single-exponential time algorithm for computing nonnegative rank*, `arXiv:1205.0044`.

[19] D. Mond, J. Smith, and D. van Straten: Stochastic factorizations, sandwiched simplices and the topology of the space of explanations, *R. Soc. Lond. Proc. Ser. A Math. Phys. Eng. Sci.* **459** (2003) 2821–2845.

[20] A.P. Morgan and A.J. Sommese: *A homotopy for solving general polynomial systems that respects m-homogeneous structures*, Appl. Math. Comput. **24** (1987) 101–113.

[21] L. Pachter and B. Sturmfels: *Algebraic Statistics for Computational Biology*, Cambridge University Press, 2005.

[22] A.J. Sommese and C.W. Wampler: *The Numerical Solution of Systems of Polynomials Arising in Engineering and Science*, World Scientific, Singapore, 2005.

[23] A.J. Sommese, J. Verschelde, and C.W. Wampler: *Symmetric functions applied to decomposing solution sets of polynomial systems*, SIAM J. Numer. Anal. **40** (2002) 2026–2046.

[24] S. Steidel: Gröbner bases of symmetric ideals, Journal of Symbolic Computation **54** (2013) 72–86.

[25] W. Stein et al: *Sage Mathematics Software* (Version 5.0), The Sage Development Team, 2012, `http://www.sagemath.org`.

[26] B. Sturmfels: *Open problems in algebraic statistics*, in "Emerging Applications of Algebraic Geometry", (editors M. Putinar and S. Sullivant), I.M.A. Volumes in Mathematics and its Applications, 149, Springer, New York, 2008, pp. 351–364.

[27] J. Verschelde: *Algorithm 795: PHCpack: a general-purpose solver for polynomial systems by homotopy continuation*, ACM Trans. Math. Software **25** (1999) 251–276.

[28] J. Verschelde and R. Cools: *Symbolic homotopy construction*, Appl. Algebra Engrg. Comm. Comput. **4** (1993) 169–183.

[29] M. Zhu, G. Jiang and S. Gao: *Solving the 100 Swiss Francs problem*, Mathematics in Computer Science **5** (2011) 195–207.