Numerical Multilinear Algebra I

Lek-Heng Lim

University of California, Berkeley

January 5-7, 2009

L.-H. Lim (ICM Lecture)

Numerical Multilinear Algebra I

January 5-7, 2009 1 / 55

A 🖓 h

3

Hope

Past 50 years, Numerical Linear Algebra played indispensable role in

- the statistical analysis of two-way data,
- the numerical solution of partial differential equations arising from vector fields,
- the numerical solution of second-order optimization methods.

Next step — development of Numerical Multilinear Algebra for

- the statistical analysis of multi-way data,
- the numerical solution of partial differential equations arising from tensor fields,
- the numerical solution of higher-order optimization methods.

- 31

くほと くほと くほと

DARPA mathematical challenge eight

One of the twenty three mathematical challenges announced at DARPA Tech 2007.

Problem

Beyond convex optimization: *can linear algebra be replaced by algebraic geometry in a systematic way?*

- Algebraic geometry in a slogan: polynomials are to algebraic geometry what matrices are to linear algebra.
- Polynomial $f \in \mathbb{R}[x_1, \dots, x_n]$ of degree d can be expressed as

$$f(\mathbf{x}) = a_0 + \mathbf{a}_1^\top \mathbf{x} + \mathbf{x}^\top A_2 \mathbf{x} + A_3(\mathbf{x}, \mathbf{x}, \mathbf{x}) + \cdots + A_d(\mathbf{x}, \dots, \mathbf{x}).$$

 $a_0 \in \mathbb{R}, a_1 \in \mathbb{R}^n, A_2 \in \mathbb{R}^{n \times n}, A_3 \in \mathbb{R}^{n \times n \times n}, \dots, A_d \in \mathbb{R}^{n \times \dots \times n}.$

- Numerical linear algebra: d = 2.
- Numerical multilinear algebra: d > 2.

Motivation

Why multilinear:

- "Classification of mathematical problems as linear and nonlinear is like classification of the Universe as bananas and non-bananas."
- Nonlinear too general. Multilinear next natural step.

Why numerical:

- Different from Computer Algebra.
- Numerical rather than symbolic: floating point operations cheap and abundant; symbolic operations expensive.
- Like other areas in numerical analysis, will entail the approximate solution of approximate multilinear problems with approximate data but under controllable and rigorous confidence bounds on the errors involved.

- 3

Tensors: mathematician's definition

 U, V, W vector spaces. Think of U ⊗ V ⊗ W as the vector space of all formal linear combinations of terms of the form u ⊗ v ⊗ w,

$$\sum \alpha \mathbf{u} \otimes \mathbf{v} \otimes \mathbf{w},$$

where $\alpha \in \mathbb{R}, \mathbf{u} \in U, \mathbf{v} \in V, \mathbf{w} \in W$.

• One condition: \otimes decreed to have the multilinear property

$$(\alpha \mathbf{u}_1 + \beta \mathbf{u}_2) \otimes \mathbf{v} \otimes \mathbf{w} = \alpha \mathbf{u}_1 \otimes \mathbf{v} \otimes \mathbf{w} + \beta \mathbf{u}_2 \otimes \mathbf{v} \otimes \mathbf{w},$$
$$\mathbf{u} \otimes (\alpha \mathbf{v}_1 + \beta \mathbf{v}_2) \otimes \mathbf{w} = \alpha \mathbf{u} \otimes \mathbf{v}_1 \otimes \mathbf{w} + \beta \mathbf{u} \otimes \mathbf{v}_2 \otimes \mathbf{w},$$
$$\mathbf{u} \otimes \mathbf{v} \otimes (\alpha \mathbf{w}_1 + \beta \mathbf{w}_2) = \alpha \mathbf{u} \otimes \mathbf{v} \otimes \mathbf{w}_1 + \beta \mathbf{u} \otimes \mathbf{v} \otimes \mathbf{w}_2.$$

Up to a choice of bases on U, V, W, A ∈ U ⊗ V ⊗ W can be represented by a 3-hypermatrix A = [[a_{ijk}]]^{l,m,n}_{i,j,k=1} ∈ ℝ^{l×m×n}.

Tensors: physicist's definition

- "What are tensors?" \equiv "What kind of physical quantities can be represented by tensors?"
- Usual answer: if they satisfy some 'transformation rules' under a change-of-coordinates.

Theorem (Change-of-basis)

Two representations A, A' of **A** in different bases are related by

$$(L, M, N) \cdot A = A'$$

with L, M, N respective change-of-basis matrices (non-singular).

 Pitfall: tensor fields (roughly, tensor-valued functions on manifolds) often referred to as tensors — stress tensor, piezoelectric tensor, moment-of-inertia tensor, gravitational field tensor, metric tensor, curvature tensor.

L.-H. Lim (ICM Lecture)

Tensors: data analyst's definition

- Data structure: k-array $A = \llbracket a_{ijk} \rrbracket_{i,j,k=1}^{l,m,n} \in \mathbb{R}^{l \times m \times n}$
- Algebraic structure:
 - **4** Addition/scalar multiplication: for $[\![b_{ijk}]\!] \in \mathbb{R}^{l \times m \times n}$, $\lambda \in \mathbb{R}$,

 $\llbracket a_{ijk} \rrbracket + \llbracket b_{ijk} \rrbracket := \llbracket a_{ijk} + b_{ijk} \rrbracket \quad \text{and} \quad \lambda \llbracket a_{ijk} \rrbracket := \llbracket \lambda a_{ijk} \rrbracket \in \mathbb{R}^{l \times m \times n}$

2 Multilinear matrix multiplication: for matrices $L = [\lambda_{i'i}] \in \mathbb{R}^{p \times l}, M = [\mu_{j'j}] \in \mathbb{R}^{q \times m}, N = [\nu_{k'k}] \in \mathbb{R}^{r \times n},$

$$(L, M, N) \cdot A := \llbracket c_{i'j'k'} \rrbracket \in \mathbb{R}^{p \times q \times r}$$

where

$$c_{i'j'k'} := \sum_{i=1}^{l} \sum_{j=1}^{m} \sum_{k=1}^{n} \lambda_{i'i} \mu_{j'j} \nu_{k'k} a_{ijk}.$$

- Think of A as 3-dimensional **hypermatrix**. (L, M, N) · A as multiplication on '3 sides' by matrices L, M, N.
- Generalizes to arbitrary order k. If k = 2, ie. matrix, then $(M, N) \cdot A = MAN^{T}$.

Hypermatrices

Totally ordered finite sets: $[n] = \{1 < 2 < \cdots < n\}, n \in \mathbb{N}.$

• Vector or *n*-tuple

$$f:[n] \to \mathbb{R}.$$

If $f(i) = a_i$, then f is represented by $\mathbf{a} = [a_1, \dots, a_n]^\top \in \mathbb{R}^n$. • Matrix

$$f:[m]\times [n]\to \mathbb{R}.$$

If $f(i,j) = a_{ij}$, then f is represented by $A = [a_{ij}]_{i,j=1}^{m,n} \in \mathbb{R}^{m \times n}$.

Hypermatrix (order 3)

$$f:[I]\times [m]\times [n]\to \mathbb{R}.$$

If $f(i, j, k) = a_{ijk}$, then f is represented by $\mathcal{A} = [\![a_{ijk}]\!]_{i,j,k=1}^{l,m,n} \in \mathbb{R}^{l \times m \times n}$. Normally $\mathbb{R}^X = \{f : X \to \mathbb{R}\}$. Ought to be $\mathbb{R}^{[n]}, \mathbb{R}^{[m] \times [n]}, \mathbb{R}^{[l] \times [m] \times [n]}$.

Hypermatrices and tensors

Up to choice of bases

- $\mathbf{a} \in \mathbb{R}^n$ can represent a vector in V (contravariant) or a linear functional in V^* (covariant).
- A ∈ ℝ^{m×n} can represent a bilinear form V* × W* → ℝ (contravariant), a bilinear form V × W → ℝ (covariant), or a linear operator V → W (mixed).
- $\mathcal{A} \in \mathbb{R}^{l \times m \times n}$ can represent trilinear form $U \times V \times W \to \mathbb{R}$ (covariant), bilinear operators $V \times W \to U$ (mixed), etc.

A hypermatrix is the same as a tensor if

- we give it coordinates (represent with respect to some bases);
- 2 we ignore covariance and contravariance.

Basic operation on a hypermatrix

• A matrix can be multiplied on the left and right: $A \in \mathbb{R}^{m \times n}$, $X \in \mathbb{R}^{p \times m}$, $Y \in \mathbb{R}^{q \times n}$,

$$(X, Y) \cdot A = XAY^{\top} = [c_{\alpha\beta}] \in \mathbb{R}^{p \times q}$$

where

$$c_{lphaeta} = \sum_{i,j=1}^{m,n} x_{lpha i} y_{eta j} \mathsf{a}_{ij}.$$

• A hypermatrix can be multiplied on three sides: $\mathcal{A} = \llbracket a_{ijk} \rrbracket \in \mathbb{R}^{l \times m \times n}$, $X \in \mathbb{R}^{p \times l}$, $Y \in \mathbb{R}^{q \times m}$, $Z \in \mathbb{R}^{r \times n}$,

$$(X, Y, Z) \cdot \mathcal{A} = \llbracket c_{\alpha\beta\gamma} \rrbracket \in \mathbb{R}^{p \times q \times r}$$

where

$$c_{\alpha\beta\gamma} = \sum_{i,j,k=1}^{I,m,n} x_{\alpha i} y_{\beta j} z_{\gamma k} a_{ijk}.$$

・ 同 ト ・ ヨ ト ・ ヨ ト

Basic operation on a hypermatrix

Covariant version:

$$\mathcal{A} \cdot (X^{\top}, Y^{\top}, Z^{\top}) := (X, Y, Z) \cdot \mathcal{A}.$$

Gives convenient notations for multilinear functionals and multilinear operators. For x ∈ ℝ^l, y ∈ ℝ^m, z ∈ ℝⁿ,

$$\mathcal{A}(\mathbf{x}, \mathbf{y}, \mathbf{z}) := \mathcal{A} \cdot (\mathbf{x}, \mathbf{y}, \mathbf{z}) = \sum_{\substack{i,j,k=1\\i,j,k=1}}^{l,m,n} a_{ijk} x_i y_j z_k,$$
$$\mathcal{A}(l, \mathbf{y}, \mathbf{z}) := \mathcal{A} \cdot (l, \mathbf{y}, \mathbf{z}) = \sum_{\substack{m,n\\j,k=1}}^{m,n} a_{ijk} y_j z_k.$$

L.-H. Lim (ICM Lecture)

January 5-7, 2009 11 / 55

3

Segre outer product

If $U = \mathbb{R}^{l}$, $V = \mathbb{R}^{m}$, $W = \mathbb{R}^{n}$, $\mathbb{R}^{l} \otimes \mathbb{R}^{m} \otimes \mathbb{R}^{n}$ may be identified with $\mathbb{R}^{l \times m \times n}$ if we define \otimes by

$$\mathbf{u}\otimes\mathbf{v}\otimes\mathbf{w}=\llbracket u_iv_jw_k\rrbracket_{i,j,k=1}^{l,m,n}.$$

A tensor $A \in \mathbb{R}^{l \times m \times n}$ is said to be **decomposable** if it can be written in the form

$$A = \mathbf{u} \otimes \mathbf{v} \otimes \mathbf{w}$$

for some $\mathbf{u} \in \mathbb{R}^{l}$, $\mathbf{v} \in \mathbb{R}^{m}$, $\mathbf{w} \in \mathbb{R}^{n}$.

The set of all decomposable tensors is known as the **Segre variety** in algebraic geometry. It is a closed set (in both the Euclidean and Zariski sense) as it can be described algebraically:

$$\mathsf{Seg}(\mathbb{R}^{l},\mathbb{R}^{m},\mathbb{R}^{n}) = \{A \in \mathbb{R}^{l \times m \times n} \mid a_{i_{1}i_{2}i_{3}}a_{j_{1}j_{2}j_{3}} = a_{k_{1}k_{2}k_{3}}a_{l_{1}l_{2}l_{3}}, \{i_{\alpha},j_{\alpha}\} = \{k_{\alpha},l_{\alpha}\}\}$$

Symmetric hypermatrices

• Cubical hypermatrix $[\![a_{ijk}]\!] \in \mathbb{R}^{n \times n \times n}$ is symmetric if

$$a_{ijk} = a_{ikj} = a_{jik} = a_{jki} = a_{kij} = a_{kji}.$$

- Invariant under all permutations $\sigma \in \mathfrak{S}_k$ on indices.
- $S^k(\mathbb{R}^n)$ denotes set of all order-k symmetric hypermatrices.

Example

Higher order derivatives of multivariate functions.

Example

Moments of a random vector $\mathbf{x} = (X_1, \dots, X_n)$:

$$m_k(\mathbf{x}) = \left[E(x_{i_1}x_{i_2}\cdots x_{i_k}) \right]_{i_1,\ldots,i_k=1}^n = \left[\int \cdots \int x_{i_1}x_{i_2}\cdots x_{i_k} \ d\mu(x_{i_1})\cdots d\mu(x_{i_k}) \right]_{i_1,\ldots,i_k=1}^n.$$

< 🗇 >

Symmetric hypermatrices

Example

Cumulants of a random vector $\mathbf{x} = (X_1, \ldots, X_n)$:

$$\kappa_k(\mathbf{x}) = \left[\sum_{A_1 \sqcup \cdots \sqcup A_p = \{i_1, \dots, i_k\}} (-1)^{p-1} (p-1)! E\left(\prod_{i \in A_1} x_i\right) \cdots E\left(\prod_{i \in A_p} x_i\right)\right]_{i_1, \dots, i_k = 1}^n.$$

For n = 1, $\kappa_k(x)$ for k = 1, 2, 3, 4 are the expectation, variance, skewness, and kurtosis.

• Important in Independent Component Analysis (ICA).

くほと くほと くほと

Inner products and norms

•
$$\ell^2([n])$$
: $\mathbf{a}, \mathbf{b} \in \mathbb{R}^n$, $\langle \mathbf{a}, \mathbf{b} \rangle = \mathbf{a}^\top \mathbf{b} = \sum_{i=1}^n a_i b_i$.

• $\ell^2([m] \times [n])$: $A, B \in \mathbb{R}^{m \times n}$, $\langle A, B \rangle = tr(A^\top B) = \sum_{i,j=1}^{m,n} a_{ij} b_{ij}$.

• $\ell^2([l] \times [m] \times [n])$: $\mathcal{A}, \mathcal{B} \in \mathbb{R}^{l \times m \times n}, \langle \mathcal{A}, \mathcal{B} \rangle = \sum_{i,j,k=1}^{l,m,n} a_{ijk} b_{ijk}.$

In general,

$$\ell^{2}([m] \times [n]) = \ell^{2}([m]) \otimes \ell^{2}([n]),$$

$$\ell^{2}([l] \times [m] \times [n]) = \ell^{2}([l]) \otimes \ell^{2}([m]) \otimes \ell^{2}([n]).$$

Frobenius norm

$$\|\mathcal{A}\|_{F}^{2} = \sum_{i,j,k=1}^{l,m,n} a_{ijk}^{2}.$$

• Norm topology often more directly relevant to engineering applications than Zariski toplogy.

Other norms

Let ||·||_{αi} be a norm on ℝ^{di}, i = 1,..., k. Then operator norm of multilinear functional A : ℝ^{d1} × ··· × ℝ^{dk} → ℝ is

$$\|A\|_{\alpha_1,\ldots,\alpha_k} := \sup \frac{|A(\mathbf{x}_1,\ldots,\mathbf{x}_k)|}{\|\mathbf{x}_1\|_{\alpha_1}\cdots\|\mathbf{x}_k\|_{\alpha_k}}.$$

Deep and important results about such norms in functional analysis. *E-norm* and *G-norm*:

$$\|A\|_E = \sum_{i_1,...,i_k=1}^{d_1,...,d_k} |a_{j_1\cdots j_k}|$$

and

$$\|A\|_G = \max\{|a_{j_1\cdots j_k}| \mid j_1 = 1, \dots, d_1; \dots; j_k = 1, \dots, d_k\}.$$

• Multiplicative on rank-1 tensors:

$$\|\mathbf{u} \otimes \mathbf{v} \otimes \cdots \otimes \mathbf{z}\|_{E} = \|\mathbf{u}\|_{1} \|\mathbf{v}\|_{1} \cdots \|\mathbf{z}\|_{1},$$

$$\|\mathbf{u} \otimes \mathbf{v} \otimes \cdots \otimes \mathbf{z}\|_{F} = \|\mathbf{u}\|_{2} \|\mathbf{v}\|_{2} \cdots \|\mathbf{z}\|_{2},$$

$$\|\mathbf{u} \otimes \mathbf{v} \otimes \cdots \otimes \mathbf{z}\|_{G} = \|\mathbf{u}\|_{\infty} \|\mathbf{v}\|_{\infty} \cdots \|\mathbf{z}\|_{\infty}.$$

L.-H. Lim (ICM Lecture)

Tensor ranks (Hitchcock, 1927)

• Matrix rank.
$$A \in \mathbb{R}^{m \times n}$$

$$\begin{aligned} \operatorname{rank}(A) &= \operatorname{dim}(\operatorname{span}_{\mathbb{R}}\{A_{\bullet 1}, \dots, A_{\bullet n}\}) & (\operatorname{column rank}) \\ &= \operatorname{dim}(\operatorname{span}_{\mathbb{R}}\{A_{1\bullet}, \dots, A_{m\bullet}\}) & (\operatorname{row rank}) \\ &= \min\{r \mid A = \sum_{i=1}^{r} \mathbf{u}_{i} \mathbf{v}_{i}^{\top}\} & (\operatorname{outer product rank}). \end{aligned}$$

• Multilinear rank. $\mathcal{A} \in \mathbb{R}^{l \times m \times n}$. rank_{\boxplus}(\mathcal{A}) = ($r_1(\mathcal{A}), r_2(\mathcal{A}), r_3(\mathcal{A})$),

$$\begin{split} r_1(\mathcal{A}) &= \dim(\operatorname{span}_{\mathbb{R}}\{\mathcal{A}_{1\bullet\bullet}, \dots, \mathcal{A}_{I\bullet\bullet}\})\\ r_2(\mathcal{A}) &= \dim(\operatorname{span}_{\mathbb{R}}\{\mathcal{A}_{\bullet1\bullet}, \dots, \mathcal{A}_{\bulletm\bullet}\})\\ r_3(\mathcal{A}) &= \dim(\operatorname{span}_{\mathbb{R}}\{\mathcal{A}_{\bullet\bullet1}, \dots, \mathcal{A}_{\bullet\bulletn}\}) \end{split}$$

• Outer product rank. $\mathcal{A} \in \mathbb{R}^{l \times m \times n}$.

$$\operatorname{rank}_{\otimes}(\mathcal{A}) = \min\{r \mid \mathcal{A} = \sum_{i=1}^{r} \mathbf{u}_i \otimes \mathbf{v}_i \otimes \mathbf{w}_i\}$$

where $\mathbf{u} \otimes \mathbf{v} \otimes \mathbf{w} := \llbracket u_i v_j w_k \rrbracket_{i,j,k=1}^{l,m,n}$.

Properties of matrix rank

- **Q** Rank of $A \in \mathbb{R}^{m \times n}$ easy to determine (Gaussian elimination)
- ② Best rank-r approximation to A ∈ ℝ^{m×n} always exist (Eckart-Young theorem)
- Solution Best rank-*r* approximation to $A \in \mathbb{R}^{m \times n}$ easy to find (singular value decomposition)
- Pick A ∈ ℝ^{m×n} at random, then A has full rank with probability 1, ie. rank(A) = min{m, n}
- rank(A) from a non-orthogonal rank-revealing decomposition (e.g. A = L₁DL₂^T) and rank(A) from an orthogonal rank-revealing decomposition (e.g. A = Q₁RQ₂^T) are equal
- o rank(A) is base field independent, i.e. same value whether we regard A as an element of ℝ^{m×n} or as an element of ℂ^{m×n}

Properties of outer product rank

- **(**) Computing rank_{\otimes}(*A*) for *A* \in $\mathbb{R}^{l \times m \times n}$ is **NP-hard** [Håstad 1990]
- Solution
 Solution
- When $\operatorname{argmin}_{\operatorname{rank}_{\otimes}(B) \leq r} ||A B||_F$ does have a solution, computing the solution is an **NP-complete** problem in general
- For some *l*, *m*, *n*, if we sample A ∈ ℝ^{l×m×n} at random, there is no r such that rank_⊗(A) = r with probability 1
- Solution A constraints on X, Y, Z will in general require a sum with more than rank_⊗(A) number of terms
- In rank_⊗(A) is **base field dependent**, i.e. value depends on whether we regard A ∈ ℝ^{l×m×n} or A ∈ ℂ^{l×m×n}

▲□▶ ▲□▶ ▲□▶ ▲□▶ □ のの⊙

Properties of multilinear rank

- Computing rank_{\boxplus}(A) for $A \in \mathbb{R}^{l \times m \times n}$ is easy
- Solution to $\operatorname{argmin}_{\operatorname{rank}_{\boxplus}(B) \leq (r_1, r_2, r_3)} ||A B||_F$ always exist
- Solution to $\operatorname{argmin}_{\operatorname{rank}_{\mathbb{H}}(B) \leq (r_1, r_2, r_3)} ||A B||_F$ easy to find
- Pick $A \in \mathbb{R}^{l \times m \times n}$ at random, then A has

 $\operatorname{rank}_{\boxplus}(A) = (\min(I, mn), \min(m, ln), \min(n, lm))$

with probability 1

- If A ∈ ℝ^{l×m×n} has rank_⊞(A) = (r₁, r₂, r₃). Then there exist full-rank matrices X ∈ ℝ^{l×r₁}, Y ∈ ℝ^{m×r₂}, Z ∈ ℝ^{n×r₃} and core tensor C ∈ ℝ<sup>r₁×r₂×r₃ such that A = (X, Y, Z) · C. X, Y, Z may be chosen to have orthonormal columns
 </sup>
- I rank_⊞(A) is base field independent, i.e. same value whether we regard A ∈ ℝ^{l×m×n} or A ∈ ℂ^{l×m×n}

▲□▶ ▲□▶ ▲□▶ ▲□▶ □ のの⊙

Algebraic computational complexity

• For
$$A = (a_{ij}), B = (b_{jk}) \in \mathbb{R}^{n \times n},$$

 $AB = \sum_{i,j,k=1}^{n} a_{ik} b_{kj} E_{ij} = \sum_{i,j,k=1}^{n} \varphi_{ik}(A) \varphi_{kj}(B) E_{ij}$
where $E_{ij} = \mathbf{e}_i \mathbf{e}_j^\top \in \mathbb{R}^{n \times n}.$ Let
 $T = \sum_{i,j,k=1}^{n} \varphi_{ik} \otimes \varphi_{kj} \otimes E_{ij}.$

- $O(n^{2+\varepsilon})$ algorithm for multiplying two $n \times n$ matrices gives $O(n^{2+\varepsilon})$ algorithm for solving system of n linear equations [Strassen 1969].
- Conjecture. $\log_2(\operatorname{rank}_{\otimes}(T)) \leq 2 + \varepsilon$.
- Best known result. $O(n^{2.376})$ [Coppersmith-Winograd 1987; Cohn-Kleinberg-Szegedy-Umans 2005].

More tensor ranks

• For
$$\mathbf{u} \in \mathbb{R}^{l}, \mathbf{v} \in \mathbb{R}^{m}, \mathbf{w} \in \mathbb{R}^{n}$$
,

$$\mathbf{u} \otimes \mathbf{v} \otimes \mathbf{w} := \llbracket u_i v_j w_k \rrbracket_{i,j,k=1}^{l,m,n} \in \mathbb{R}^{l \times m \times n}.$$

• Outer product rank. $\mathcal{A} \in \mathbb{R}^{l \times m \times n}$,

$$\operatorname{rank}_{\otimes}(\mathcal{A}) = \min\{r \mid \mathcal{A} = \sum_{i=1}^{r} \sigma_i \mathbf{u}_i \otimes \mathbf{v}_i \otimes \mathbf{w}_i, \quad \sigma_i \in \mathbb{R}\}.$$

• Symmetric outer product rank. $\mathcal{A} \in \mathsf{S}^k(\mathbb{R}^n)$,

$$\mathsf{rank}_{\mathsf{S}}(\mathcal{A}) = \min\{r \mid \mathcal{A} = \sum_{i=1}^{r} \lambda_i \mathbf{v}_i \otimes \mathbf{v}_i \otimes \mathbf{v}_i, \quad \lambda_i \in \mathbb{R}\}.$$

• Nonnegative outer product rank. $\mathcal{A} \in \mathbb{R}^{l imes m imes n}_+$,

$$\mathsf{rank}_+(\mathcal{A}) = \min\{r \mid \mathcal{A} = \sum_{i=1}^r \delta_i \mathbf{x}_i \otimes \mathbf{y}_i \otimes \mathbf{z}_i, \quad \delta_i \in \mathbb{R}_+\}.$$

・ 同 ト ・ ヨ ト ・ ヨ ト … ヨ …

SVD, EVD, NMF of a matrix

• Singular value decomposition of $A \in \mathbb{R}^{m \times n}$,

$$A = U \Sigma V^{\top} = \sum_{i=1}^{r} \sigma_i \mathbf{u}_i \otimes \mathbf{v}_i$$

where rank(A) = r, $U \in O(m)$ left singular vectors, $V \in O(n)$ right singular vectors, Σ singular values.

• Symmetric eigenvalue decomposition of $A \in S^2(\mathbb{R}^n)$,

$$A = V \Lambda V^{\top} = \sum_{i=1}^{r} \lambda_i \mathbf{v}_i \otimes \mathbf{v}_i,$$

where rank(A) = r, $V \in O(n)$ eigenvectors, Λ eigenvalues.

• Nonnegative matrix factorization of $A \in \mathbb{R}^{n \times n}_+$,

$$A = X \Delta Y^{\top} = \sum_{i=1}^{r} \delta_i \mathbf{x}_i \otimes \mathbf{y}_i$$

where rank₊(A) = r, $X, Y \in \mathbb{R}^{m \times r}_+$ unit column vectors (in the 1-norm), Δ positive values.

L.-H. Lim (ICM Lecture)

SVD, EVD, NMF of a hypermatrix

• Outer product decomposition of $\mathcal{A} \in \mathbb{R}^{l \times m \times n}$,

$$\mathcal{A} = \sum_{i=1}^r \sigma_i \mathbf{u}_i \otimes \mathbf{v}_i \otimes \mathbf{w}_i$$

where $\operatorname{rank}_{\otimes}(\mathcal{A}) = r$, $\mathbf{u}_i \in \mathbb{R}^l$, $\mathbf{v}_i \in \mathbb{R}^m$, $\mathbf{w}_i \in \mathbb{R}^n$ unit vectors, $\sigma_i \in \mathbb{R}$.

• Symmetric outer product decomposition of $\mathcal{A} \in S^3(\mathbb{R}^n)$,

$$\mathcal{A} = \sum_{i=1}^r \lambda_i \mathbf{v}_i \otimes \mathbf{v}_i \otimes \mathbf{v}_i$$

where rank_S(A) = r, \mathbf{v}_i unit vector, $\lambda_i \in \mathbb{R}$.

• Nonnegative outer product decomposition for hypermatrix $\mathcal{A} \in \mathbb{R}^{l \times m \times n}_+$ is

$$\mathcal{A} = \sum_{i=1}^r \delta_i \mathbf{x}_i \otimes \mathbf{y}_i \otimes \mathbf{z}_i$$

where rank₊(A) = r, $\mathbf{x}_i \in \mathbb{R}^l_+, \mathbf{y}_i \in \mathbb{R}^m_+, \mathbf{z}_i \in \mathbb{R}^n_+$ unit vectors, $\delta_i \in \mathbb{R}_+$. Best low rank approximation of a matrix

• Given $A \in \mathbb{R}^{m \times n}$. Want

 $\operatorname{argmin}_{\operatorname{rank}(B)\leq r} \|A - B\|.$

• More precisely, find σ_i , \mathbf{u}_i , \mathbf{v}_i , $i = 1, \ldots, r$, that minimizes

$$\|\mathcal{A} - \sigma_1 \mathbf{u}_1 \otimes \mathbf{v}_1 - \sigma_2 \mathbf{u}_2 \otimes \mathbf{v}_2 - \cdots - \sigma_r \mathbf{u}_r \otimes \mathbf{v}_r\|.$$

Theorem (Eckart–Young)

Let $A = U\Sigma V^{\top} = \sum_{i=1}^{\operatorname{rank}(A)} \sigma_i \mathbf{u}_i \mathbf{v}_i^{\top}$ be singular value decomposition. For $r \leq \operatorname{rank}(A)$, let

$$A_r := \sum_{i=1}^r \sigma_i \mathbf{u}_i \mathbf{v}_i^\top.$$

Then

$$\|A - A_r\|_F = \min_{\mathsf{rank}(B) \le r} \|A - B\|_F.$$

• No such thing for hypermatrices of order 3 or higher.

L.-H. Lim (ICM Lecture)

Segre variety and its secant varieties

- The set of all rank-1 hypermatrices is known as the Segre variety in algebraic geometry.
- It is a closed set (in both the Euclidean and Zariski sense) as it can be described algebraically:

$$Seg(\mathbb{R}^{l}, \mathbb{R}^{m}, \mathbb{R}^{n}) = \{ \mathcal{A} \in \mathbb{R}^{l \times m \times n} \mid \mathcal{A} = \mathbf{u} \otimes \mathbf{v} \otimes \mathbf{w} \} = \{ \mathcal{A} \in \mathbb{R}^{l \times m \times n} \mid \mathbf{a}_{i_{1}i_{2}i_{3}}\mathbf{a}_{j_{1}j_{2}j_{3}} = \mathbf{a}_{k_{1}k_{2}k_{3}}\mathbf{a}_{l_{1}l_{2}l_{3}}, \{i_{\alpha}, j_{\alpha}\} = \{k_{\alpha}, l_{\alpha}\} \}$$

- Hypermatrices that have rank > 1 are elements on the higher secant varieties of 𝒴 = Seg(ℝ^l, ℝ^m, ℝⁿ).
- E.g. a hypermatrix has rank 2 if it sits on a secant line through two points in \mathscr{S} but not on \mathscr{S} , rank 3 if it sits on a secant plane through three points in \mathscr{S} but not on any secant lines, etc.
- Minor technicality: should really be secant quasiprojective variety.

イロト 人間ト イヨト イヨト

Scientific data mining

- **Spectroscopy:** measure light absorption/emission of specimen as function of energy.
- Typical **specimen** contains 10¹³ to 10¹⁶ light absorbing entities or **chromophores** (molecules, amino acids, etc).

Fact (Beer's Law)

 $A(\lambda) = -\log(I_1/I_0) = \varepsilon(\lambda)c$. A = absorbance, $I_1/I_0 = fraction of$ intensity of light of wavelength λ that passes through specimen, c =concentration of chromophores.

Multiple chromophores (f = 1,...,r) and wavelengths (i = 1,...,m) and specimens/experimental conditions (j = 1,...,n),

$$A(\lambda_i, s_j) = \sum_{f=1}^r \varepsilon_f(\lambda_i) c_f(s_j).$$

• Bilinear model aka **factor analysis**: $A_{m \times n} = E_{m \times r} C_{r \times n}$ rank-revealing factorization or, in the presence of noise, low-rank approximation min $||A_{m \times n} - E_{m \times r} C_{r \times n}||$.

Modern data mining

- Text mining is the spectroscopy of documents.
- Specimens = **documents**.
- Chromophores = **terms**.
- Absorbance = inverse document frequency:

$$A(t_i) = -\log\left(\sum_j \chi(f_{ij})/n\right).$$

- Concentration = term frequency: f_{ij} .
- $\sum_{i} \chi(f_{ij})/n$ = fraction of documents containing t_i .
- A ∈ ℝ^{m×n} term-document matrix. A = QR = UΣV^T rank-revealing factorizations.
- Bilinear model aka vector space model.
- Due to Gerald Salton and colleagues: SMART (system for the mechanical analysis and retrieval of text).

・ 回 ト ・ ヨ ト ・ ヨ ト ・ ヨ

Bilinear models

- Bilinear models work on 'two-way' data:
 - ► measurements on object *i* (genomes, chemical samples, images, webpages, consumers, etc) yield a vector a_i ∈ ℝⁿ where n = number of features of *i*;
 - collection of *m* such objects, *A* = [a₁,..., a_m] may be regarded as an *m*-by-*n* matrix, e.g. gene × microarray matrices in bioinformatics, terms × documents matrices in text mining, facial images × individuals matrices in computer vision.
- Various matrix techniques may be applied to extract useful information: QR, EVD, SVD, NMF, CUR, compressed sensing techniques, etc.
- Examples: vector space model, factor analysis, principal component analysis, latent semantic indexing, PageRank, EigenFaces.
- Some problems: factor indeterminacy A = XY rank-revealing factorization not unique; unnatural for *k*-way data when k > 2.

- 3

イロト 人間ト イヨト イヨト

Ubiquity of multiway data

- Batch data: batch \times time \times variable
- Time-series analysis: time \times variable \times lag
- Computer vision: people \times view \times illumination \times expression \times pixel
- **Bioinformatics:** gene × microarray × oxidative stress
- Phylogenetics: $codon \times codon \times codon$
- Analytical chemistry: sample \times elution time \times wavelength
- Atmospheric science: location × variable × time × observation
- Psychometrics: individual × variable × time
- Sensory analysis: sample \times attribute \times judge
- Marketing: product × product × consumer

Fact (Inevitable consequence of technological advancement) Increasingly sophisticated instruments, sensor devices, data collecting and experimental methodologies lead to increasingly complex data.

L.-H. Lim (ICM Lecture)

Numerical Multilinear Algebra I

Fundamental problem of multiway data analysis

- $\mathcal A$ hypermatrix, symmetric hypermatrix, or nonnegative hypermatrix.
- Solve

$$\operatorname{argmin}_{\operatorname{rank}(\mathcal{B})\leq r} \|\mathcal{A} - \mathcal{B}\|.$$

 rank may be outer product rank, multilinear rank, symmetric rank (for symmetric hypermatrix), or nonnegative rank (nonnegative hypermatrix).

Example

Given
$$\mathcal{A} \in \mathbb{R}^{d_1 \times d_2 \times d_3}$$
, find $\mathbf{u}_i, \mathbf{v}_i, \mathbf{w}_i, i = 1, \dots, r$, that minimizes

$$\|\mathcal{A} - \mathbf{u}_1 \otimes \mathbf{v}_1 \otimes \mathbf{w}_1 - \mathbf{u}_2 \otimes \mathbf{v}_2 \otimes \mathbf{w}_2 - \dots - \mathbf{u}_r \otimes \mathbf{v}_r \otimes \mathbf{z}_r\|$$

or $C \in \mathbb{R}^{r_1 \times r_2 \times r_3}$ and $U \in \mathbb{R}^{d_1 \times r_1}, V \in \mathbb{R}^{d_2 \times r_2}, W \in \mathbb{R}^{d_3 \times r_3}$, that minimizes

 $\|\mathcal{A} - (U, V, W) \cdot \mathcal{C}\|.$

A B A B A
 A
 B
 A
 A
 B
 A
 A
 B
 A
 A
 B
 A
 A
 B
 A
 A
 B
 A
 A
 B
 A
 A
 B
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A

Fundamental problem of multiway data analysis

Example

Given $\mathcal{A} \in S^k(\mathbb{C}^n)$, find \mathbf{u}_i , $i = 1, \ldots, r$, that minimizes

$$\|\mathcal{A} - \mathbf{u}_1^{\otimes k} - \mathbf{u}_2^{\otimes k} - \cdots - \mathbf{u}_r^{\otimes k}\|$$

or $\mathcal{C} \in \mathbb{R}^{r_1 imes r_2 imes r_3}$ and $U \in \mathbb{R}^{n imes r_i}$ that minimizes

 $\|\mathcal{A} - (U, U, U) \cdot \mathcal{C}\|.$

L.-H. Lim (ICM Lecture)

・ 同 ト ・ ヨ ト ・ ヨ ト … ヨ …

Outer product decomposition in spectroscopy

- Application to fluorescence spectral analysis by [Bro; 1997].
- Specimens with a number of pure substances in different concentration
 - a_{ijk} = fluorescence emission intensity at wavelength λ_j^{em} of *i*th sample excited with light at wavelength λ_k^{ex}.
 - Get 3-way data $\mathcal{A} = \llbracket a_{ijk} \rrbracket \in \mathbb{R}^{l \times m \times n}$.
 - Get outer product decomposition of ${\cal A}$

$$\mathcal{A} = \mathbf{x}_1 \otimes \mathbf{y}_1 \otimes \mathbf{z}_1 + \dots + \mathbf{x}_r \otimes \mathbf{y}_r \otimes \mathbf{z}_r.$$

- Get the true chemical factors responsible for the data.
 - r: number of pure substances in the mixtures,
 - x_α = (x_{1α},..., x_{lα}): relative concentrations of αth substance in specimens 1,..., l,
 - $\mathbf{y}_{\alpha} = (y_{1\alpha}, \dots, y_{m\alpha})$: excitation spectrum of α th substance,
 - $\mathbf{z}_{\alpha} = (z_{1\alpha}, \dots, z_{n\alpha})$: emission spectrum of α th substance.

• Noisy case: find best rank-*r* approximation (CANDECOMP/PARAFAC).

Uniqueness of tensor decompositions

M ∈ ℝ^{m×n}, spark(M) = size of minimal linearly dependent subset of column vectors [Donoho, Elad; 2003].

Theorem (Kruskal)

 $X = [\mathbf{x}_1, \dots, \mathbf{x}_r], Y = [\mathbf{y}_1, \dots, \mathbf{y}_r], Z = [\mathbf{z}_1, \dots, \mathbf{z}_r].$ Decomposition is unique up to scaling if

$$\operatorname{spark}(X) + \operatorname{spark}(Y) + \operatorname{spark}(Z) \ge 2r + 5.$$

- May be generalized to arbitrary order [Sidiroupoulos, Bro; 2000].
- Avoids factor indeterminacy under mild conditions.

(人間) とうき くうとう う

Multilinear decomposition in bioinformatics

- Application to cell cycle studies [Omberg, Golub, Alter; 2008].
- Collection of gene-by-microarray matrices $A_1, \ldots, A_l \in \mathbb{R}^{m \times n}$ obtained under varying oxidative stress.
 - a_{ijk} = expression level of *j*th gene in *k*th microarray under *i*th stress.
 - Get 3-way data array $\mathcal{A} = \llbracket a_{ijk} \rrbracket \in \mathbb{R}^{l \times m \times n}$.
 - Get multilinear decomposition of ${\cal A}$

$$\mathcal{A} = (X, Y, Z) \cdot \mathcal{C},$$

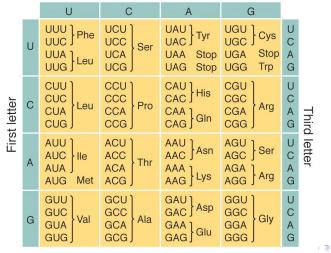
to get orthogonal matrices X, Y, Z and core tensor C by applying SVD to various 'flattenings' of A.

- Column vectors of X, Y, Z are 'principal components' or 'parameterizing factors' of the spaces of stress, genes, and microarrays; C governs interactions between these factors.
- Noisy case: approximate by discarding small c_{ijk} (Tucker Model).

イロト 不得下 イヨト イヨト 二日

Code of life is a 3-tensor

- **Codons:** triplets of nucleotides, (i, j, k) where $i, j, k \in \{A, C, G, U\}$.
- **Genetic code:** these $4^3 = 64$ codons encode the 20 amino acids.



Second letter

L.-H. Lim (ICM Lecture)

January 5-7, 2009 36 / 55

Tensors in algebraic statistical biology

Problem (Salmon conjecture)

Find the polynomial equations that defines the set

$$\{P \in \mathbb{C}^{4 \times 4 \times 4} \mid \operatorname{rank}_{\otimes}(P) \leq 4\}.$$

 Why interested? Here P = [[p_{ijk}]] is understood to mean 'complexified' probability density values with i, j, k ∈ {A, C, G, T} and we want to study tensors that are of the form

$$P = \rho_A \otimes \sigma_A \otimes \theta_A + \rho_C \otimes \sigma_C \otimes \theta_C + \rho_G \otimes \sigma_G \otimes \theta_G + \rho_T \otimes \sigma_T \otimes \theta_T,$$

in other words,

$$p_{ijk} = \rho_{Ai}\sigma_{Aj}\theta_{Ak} + \rho_{Ci}\sigma_{Cj}\theta_{Ck} + \rho_{Gi}\sigma_{Gj}\theta_{Gk} + \rho_{Ti}\sigma_{Tj}\theta_{Tk}.$$

- Why over \mathbb{C} ? Easier to deal with mathematically.
- Ultimately, want to study this over \mathbb{R}_+ .

L.-H. Lim (ICM Lecture)