

TENSOR METHODS FOR ML

- COMPRESS NEURAL NETWORKS
- LEARNING LATENT VARIABLE MODELS

I EFFICIENT COMPUTATIONS WITH TENSOR NETWORKS

- Inner product: $u, v \in \mathbb{R}^{d^N}$

What is the computational cost for $\langle u, v \rangle = \sum_{i=1}^{d^N} u_i v_i$?
 $\hookrightarrow \mathcal{O}(d^N)$

Suppose u and v are in the TT format:

$$u, v \in \mathbb{R}^{d^N} \cong \mathbb{R}^{\underbrace{d \times d \times \dots \times d}_{N \text{ times}}}$$

$$\begin{aligned} u &= \underbrace{U_1 \xrightarrow{R} U_2 \xrightarrow{R} U_3 \xrightarrow{R} \dots \xrightarrow{R} U_N}_{\downarrow \quad \downarrow \quad \downarrow \quad \downarrow} \\ v &= \underbrace{V_1 \xrightarrow{R} V_2 \xrightarrow{R} V_3 \xrightarrow{R} \dots \xrightarrow{R} V_N}_{\downarrow \quad \downarrow \quad \downarrow \quad \downarrow} \end{aligned} \quad \left. \right\} \text{TT representations of } u \text{ and } v$$

$$\langle u, v \rangle = u - v = \underbrace{U_1 \xrightarrow{R} U_2 \xrightarrow{R} U_3 \xrightarrow{R} \dots \xrightarrow{R} U_N}_{\downarrow \quad \downarrow \quad \downarrow \quad \downarrow} \quad \underbrace{V_1 \xrightarrow{R} V_2 \xrightarrow{R} V_3 \xrightarrow{R} \dots \xrightarrow{R} V_N}_{\downarrow \quad \downarrow \quad \downarrow \quad \downarrow}$$

Recall: Computing the matrix product $A B$ has cost $\mathcal{O}(mnk)$

$$(d^2 \times R) \quad (R \times d^2) \quad mn \times k$$

For ex, Computing $\underbrace{A \xrightarrow{R}}_{d \times 1}, \underbrace{B \xrightarrow{1 \times d}}_{1 \times d} \rightsquigarrow \underbrace{A \xrightarrow{R} B \xrightarrow{1 \times d}}_{d \times 1}$ can be done

with a matrix product between matrices of shapes $d^2 \times R$ and $R \times d^2$. Hence the complexity is $\mathcal{O}(Rd^4)$.

We want to compute

$$U_1 \xrightarrow{R} U_2 \xrightarrow{R} U_3 \xrightarrow{R} \dots \xrightarrow{R} U_N$$

$$V_1 \xrightarrow{R} V_2 \xrightarrow{R} V_3 \xrightarrow{R} \dots \xrightarrow{R} V_N$$

COMPLEXITY

① Compute $U_1 \xrightarrow{R}$

$$\begin{array}{c} d \\ | \\ U_1 \xrightarrow{R} \end{array}$$

$\mathcal{O}(dR^2)$

$$\begin{array}{c} R \times R \\ | \\ U_1 \xrightarrow{R} \end{array} \quad \begin{array}{c} R \times dR \\ | \\ U_2 \xrightarrow{R} \end{array}$$

$$\begin{array}{c} d \\ | \\ V_1 \xrightarrow{R} \end{array}$$

$\mathcal{O}(dR^3)$

③ $U_1 \xrightarrow{R} U_2 \xrightarrow{R}$

$$\begin{array}{c} d \\ | \\ U_1 \xrightarrow{R} \end{array} \quad \begin{array}{c} d \\ | \\ U_2 \xrightarrow{R} \end{array}$$

$$\begin{array}{c} V_1 \xrightarrow{R} \end{array} \quad \begin{array}{c} V_2 \xrightarrow{R} \\ | \\ R \times dR \end{array}$$

$\mathcal{O}(dR^3)$

$$\begin{array}{c} R \times R \\ | \\ U_1 \xrightarrow{R} U_2 \xrightarrow{R} \end{array} \quad \begin{array}{c} R \times dR \\ | \\ U_3 \xrightarrow{R} \end{array}$$

$$\begin{array}{c} d \\ | \\ V_1 \xrightarrow{R} V_2 \xrightarrow{R} \end{array} \quad \begin{array}{c} R \times dR \\ | \\ R \times dR \end{array}$$

$\mathcal{O}(dR^3)$

⑤

:

:

$$U_1 \xrightarrow{R} U_2 \xrightarrow{R} U_3 \xrightarrow{R} \dots \xrightarrow{R} U_N$$

$$V_1 \xrightarrow{R} V_2 \xrightarrow{R} V_3 \xrightarrow{R} \dots \xrightarrow{R} V_N$$

this is linear in N

TOTAL COMPLEXITY

$\mathcal{O}(NdR^3)$

↳ in contrast with $\mathcal{O}(d^N)$!
 ↗ exp. in N

- Matrix vector product
 $A \in \mathbb{R}^{d^N \times n^N}$, $x \in \mathbb{R}^{n^N}$: we want to compute $b = Ax \in \mathbb{R}^{d^N}$.

In fact, we want to compute b in the TT format.

+ TT representation of matrices :

$$A \in \mathbb{R}^{d^N \times n^N} \approx \mathbb{R}^{\underbrace{d \times d \times \dots \times d}_{N \text{ times}} \times \underbrace{n \times n \times \dots \times n}_{N \text{ times}}}$$

① First solution :

$$A = A_1 \frac{R}{d} A_2 \frac{R}{d} \dots \frac{R}{d} A_N \frac{R}{n} A_{N+1} \frac{R}{n} A_{N+2} \frac{R}{n} \dots \frac{R}{n} A_{2N}$$

↳ PBM : A will be of rank at most R .

$$A = \begin{matrix} A_1 \frac{R}{d} A_2 \frac{R}{d} \dots \frac{R}{d} A_N \end{matrix} \frac{R}{n} \begin{matrix} A_{N+1} \frac{R}{n} A_{N+2} \frac{R}{n} \dots \frac{R}{n} A_{2N} \end{matrix}$$

P Q

$d^N \times R$ $R \times n^N$

$$\hookrightarrow A = \begin{matrix} d^N \times n^N \\ / \quad \backslash \end{matrix} P Q \begin{matrix} R \times n^N \end{matrix} \rightarrow \text{Rank } R \text{ factorization of } A, \text{ hence } \text{rank}(A) \leq R.$$

② TT matrix representation :

$$A = \begin{matrix} d \\ -A_1 \frac{n}{R} \\ d \\ -A_2 \frac{n}{R} \\ d \\ -A_3 \frac{n}{R} \\ \vdots \\ d \\ -A_N \frac{n}{R} \end{matrix}$$

↳ Remark • If $A = \begin{matrix} d^2 \times n^2 \\ -A_1 \frac{n}{R} \\ -A_2 \frac{n}{R} \end{matrix}$ and $R = 1$,

$$\text{then } A = \begin{matrix} d^2 \times n^2 \\ -A_1 \frac{n}{R} \\ -A_2 \frac{n}{R} \end{matrix} = A_1 \otimes A_2 \quad \text{KRONECKER PRODUCT}$$

• In a TN, an edge of dimension 1 is equivalent to having no edge :

$$\left(\begin{matrix} -B & C \\ m \times 1 & 1 \times m \end{matrix} \right)_{ij} = \sum_{k=1}^1 B_{ik} C_{kj} = B_{i1} C_{j1} = b \circ c$$

vectors b c^T

Back to matrix vector product:

$$A = \begin{array}{c} d \\ d^4 \times h^4 \end{array} = \begin{array}{c} d \\ d \\ d \\ d \end{array} \begin{array}{c} A_1 \\ A_2 \\ A_3 \\ A_4 \end{array} \begin{array}{c} h \\ h \\ h \\ h \end{array}$$

$$x = \begin{array}{c} h \\ h \\ h \\ h \end{array} \begin{array}{c} x_1 \\ x_2 \\ x_3 \\ x_4 \end{array}$$

We want a TT representation of $b = Ax$

	$d \times R \times R$	$d \times R^2$	COMPLEXITY
$b =$	$\begin{array}{c} d \\ d^4 \end{array} \begin{array}{c} A_1 \\ A_2 \\ A_3 \\ A_4 \end{array} \begin{array}{c} h \\ h \\ h \\ h \end{array} \begin{array}{c} x_1 \\ x_2 \\ x_3 \\ x_4 \end{array}$	$= \begin{array}{c} d \\ d \end{array} \begin{array}{c} B_1 \\ B_2 \end{array} \begin{array}{c} R^2 \\ R^2 \end{array}$	$O(R^2 d h)$
$b =$	$\begin{array}{c} d \\ d^4 \end{array} \begin{array}{c} A_1 \\ A_2 \\ A_3 \\ A_4 \end{array} \begin{array}{c} h \\ h \\ h \\ h \end{array} \begin{array}{c} x_1 \\ x_2 \\ x_3 \\ x_4 \end{array}$	$= \begin{array}{c} d \\ d \end{array} \begin{array}{c} B_2 \\ B_3 \end{array} \begin{array}{c} R^2 \\ R^2 \end{array}$	$O(R^4 d h)$
$b =$	$\begin{array}{c} d \\ d^4 \end{array} \begin{array}{c} B_1 \\ B_2 \\ B_3 \\ B_4 \end{array} \begin{array}{c} R^2 \\ R^2 \\ R^2 \\ R^2 \end{array} \begin{array}{c} x_1 \\ x_2 \\ x_3 \\ x_4 \end{array}$	$= \begin{array}{c} d \\ d \end{array} \begin{array}{c} B_3 \\ B_4 \end{array} \begin{array}{c} R^2 \\ R^2 \end{array}$	$O(R^4 d h)$
			$O(R^2 d h)$

\Rightarrow If $A \in \mathbb{R}^{d^N \times h^N}$ given as a TT matrix with rank R

and $x \in \mathbb{R}^{h^N}$ given as a TT vector with rank R

then a rank R^2 TT representation of $b = Ax$ can be computed in time $O(N R^4 d h)$.

↳ in contrast with $O(d^N h^N)$

- $u, v \in \mathbb{R}^{d^N}$, given as rank R TT-vector, a rank $2R$ TT representation of $u + v$ can be computed in $\mathcal{O}(NR^3d)$.
- $u, v \in \mathbb{R}^{d^N}$, given as rank R TT-vector, a rank R^2 TT representation of $u \otimes v$ can be computed in $\mathcal{O}(NR^3d)$.
↳ component-wise product.

- TT Rounding:

Given a rank R representation of $u \in \mathbb{R}^{d^N}$, the TT rounding operation returns a rank \hat{R} TT representation of u which is "close" to the original representation (in time $\mathcal{O}(NR^3d)$).

FOR MORE DETAILS:

SIAM J. SCI. COMPUT.
Vol. 33, No. 5, pp. 2295–2317

© 2011 Society for Industrial and Applied Mathematics

TENSOR-TRAIN DECOMPOSITION*

I. V. OSELEDETS†

Abstract. A simple nonrecursive form of the tensor decomposition in d dimensions is presented. It does not inherently suffer from the curse of dimensionality; it has asymptotically the same number of parameters as the canonical decomposition, but it is stable and its computation is based on low-rank approximation of auxiliary *unfolding matrices*. The new form gives a clear and convenient way to implement all basic operations efficiently. A fast rounding procedure is presented, as well as basic linear algebra operations. Examples showing the benefits of the decomposition are given, and the efficiency is demonstrated by the computation of the smallest eigenvalue of a 19-dimensional operator.

II COMPRESSING NEURAL NETWORKS

Tensorizing Neural Networks

(NeurIPS, 2015)

Alexander Novikov^{1,4} Dmitry Podoprikin¹ Anton Osokin² Dmitry Vetrov^{1,3}

¹Skolkovo Institute of Science and Technology, Moscow, Russia

²INRIA, SIERRA project-team, Paris, France

³National Research University Higher School of Economics, Moscow, Russia

⁴Institute of Numerical Mathematics of the Russian Academy of Sciences, Moscow, Russia

novikov@bayesgroup.ru podoprikin.dmitry@gmail.com

anton.osokin@inria.fr vetrov@yandex.ru

Fully connected layer:

size: $M \times N$

$$y = \sigma(Wx + b)$$

activation function

For sake of simplicity, we ignore the bias term: $y = \sigma(Wx)$

IDEA: Represent W as a TT matrix.

$$M = m_1 \times m_2 \times \dots \times m_d$$

(d is the order of the tensor)

$$N = n_1 \times n_2 \times \dots \times n_d$$

$$W = \begin{matrix} & m_1 & m_1 \\ & G_1 & | \\ m_2 & G_2 & m_2 \\ & | & | \\ & \vdots & \vdots \\ & G_d & m_d \end{matrix}$$

TT layer

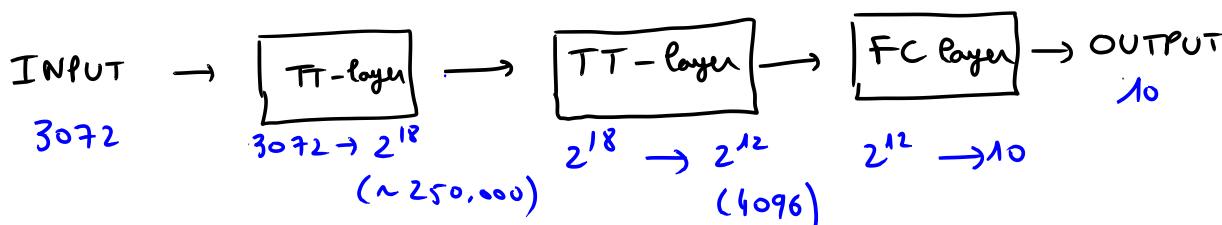
↳ parameters are
 G_1, G_2, \dots, G_d

COMPLEXITY:

Operation	Time	Memory
FC forward pass	$O(MN)$	$O(MN)$
TT forward pass	$O(dr^2 m \max\{M, N\})$	$O(r \max\{M, N\})$
FC backward pass	$O(MN)$	$O(MN)$
TT backward pass	$O(d^2 r^4 m \max\{M, N\})$	$O(r^3 \max\{M, N\})$

Table 1: Comparison of the asymptotic complexity and memory usage of an $M \times N$ TT-layer and an $M \times N$ fully-connected layer (FC). The input and output tensor shapes are $m_1 \times \dots \times m_d$ and $n_1 \times \dots \times n_d$ respectively ($m = \max_{k=1 \dots d} m_k$) and r is the maximal TT-rank.

EXPERIMENT ON CIFAR 10:



COMPUTING GRADIENTS OF TENSOR NETWORKS

↳ See TN-gradients.hdf



Tensor Decompositions for Learning Latent Variable Models

Animashree Anandkumar
*Electrical Engineering and Computer Science
 University of California, Irvine
 2200 Engineering Hall
 Irvine, CA 92697*

A.ANANDKUMAR@UCI.EDU

(JMLR, 2014)

Rong Ge
*Microsoft Research
 One Memorial Drive
 Cambridge, MA 02142*

RONGGE@MICROSOFT.COM

Daniel Hsu
*Department of Computer Science
 Columbia University
 1214 Amsterdam Avenue, #0401
 New York, NY 10027*

DJHSU@CS.COLUMBIA.EDU

Sham M. Kakade
*Microsoft Research
 One Memorial Drive
 Cambridge, MA 02142*

SKAKADE@MICROSOFT.COM

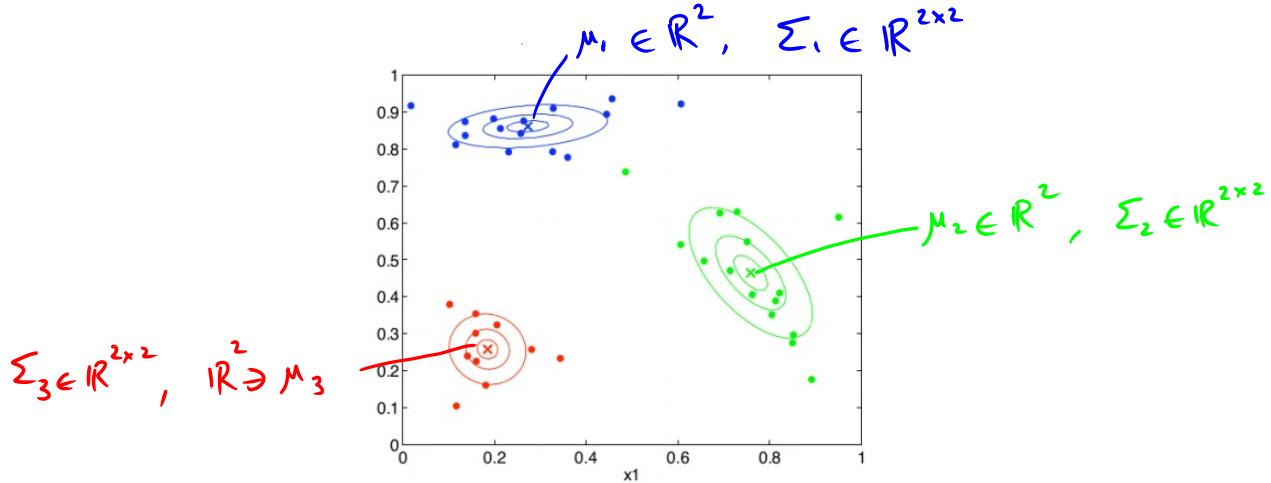
Matus Telgarsky
*Department of Statistics
 Rutgers University
 110 Frelinghuysen Road
 Piscataway, NJ 08854*

MTELGARS@CS.UCSD.EDU

Editor: Benjamin Recht

1) Examples of LVM

• Gaussian mixture model (GMM)

GMM with k components in \mathbb{R}^d :• parameters : $\mu_1, \dots, \mu_k \in \mathbb{R}^d$ (centers)n.s.d. matrices $\Sigma_1, \dots, \Sigma_k \in \mathbb{R}^{d \times d}$ (covariance matrices) $0 \leq \pi_i \leq 1, \sum_{i=1}^k \pi_i = 1 \leftarrow \pi_1, \dots, \pi_k \in \mathbb{R}_+$ (mixing probabilities)

• data generating distribution

Draw $x \sim \text{GMM}$

(i) Draw a Gaussian h with probabilities

$$\Pr(h=i) = \pi_i, i=1\dots,k$$

(ii) Draw $x \sim \mathcal{N}(\mu_h, \Sigma_h)$

• Single topic model

(latent variable is the topic of a document)

Topic: discrete random variable taking values $1, 2, \dots, k$

Vocabulary of size d (d different words).

→ We want to model the distribution over the vocabulary given a topic.

• Parameters:

$$P(\text{topic} = i) = \pi_i \quad \pi_1, \pi_2, \dots, \pi_k \in \mathbb{R} \quad (\text{topic probabilities})$$

$$0 \leq \pi_i \leq 1, \sum_{i=1}^k \pi_i = 1$$

$$P(\text{word} = j | \text{topic} = i) = (\mu_i)_j \quad \mu_1, \mu_2, \dots, \mu_k \in \mathbb{R}^d \quad (\text{voc. distrib. given a topic})$$

$$\text{for } i=1, \dots, k : \quad 0 \leq (\mu_i)_j \leq 1, \sum_{j=1}^d (\mu_i)_j = 1$$

• Data generating distribution

Draw a document of length ℓ from the single topic model.

Draw a topic h with probabilities $P(h = i) = \pi_i$	• Draw independently ℓ words with probabilities $P(\text{word } j \text{topic} = h) = (\mu_h)_j$
---	--

2) Moments and Latent variable models

Def: • If x is a random variable taking its values in \mathbb{R} , its n^{th} order moment is $\mathbb{E}[x^n]$.

• If x is a random variable taking its values in \mathbb{R}^d , its n^{th} order moment is $\mathbb{E}[x^{0:n}] = \mathbb{E}\left[\underbrace{x_0 x_0 \dots x_0}_{n \text{ times}}\right] \in \mathbb{R}^{\frac{d \times d \times \dots \times d}{n \text{ times}}}$.

↳ . 3rd moment : $\mathbb{E}[x_0 x_0 x_0]_{i,j,k} = \mathbb{E}[x_i x_j x_k]$

. 1st moment : $\mathbb{E}[x] \rightarrow \text{mean}$

. 2nd moment : $\mathbb{E}[x_0 x] = \mathbb{E}[xx^T]$

↳ if $\mathbb{E}[x] = 0$, this is the covariance matrix.

↳ $\mathbb{V}[x] = \mathbb{E}[x^2] - \mathbb{E}[x]^2$

$$\text{Cov}[x] = \mathbb{E}[xx^T] - \mathbb{E}[x]\mathbb{E}[x]^T$$

→ the method of moments give the same estimation as the maximum likelihood for a Gaussian distribution.

$$\text{If } x \sim \mathcal{N}(\mu, \sigma^2) \quad \mathbb{E}[x] = \mu \quad \mathbb{E}[x^2] = \sigma^2 + \mu^2 \quad \left| \begin{array}{l} \mu = \mathbb{E}[x] \approx \overbrace{\frac{1}{m} \sum_{i=1}^m x_i}^{\hat{\mu}} \\ \sigma^2 = \mu^2 - \mathbb{E}[x^2] \approx \hat{\mu}^2 - \overbrace{\frac{1}{m} \sum_{i=1}^m x_i^2}^{\hat{\sigma}^2} \end{array} \right.$$

Simple Topic Model

$$\text{parameters: } \Theta = \{\mu_1, \mu_2, \dots, \mu_K, \mu_{1,1}, \dots, \mu_{K,K}\}$$

$\downarrow \mathbb{R}$ $\downarrow \mathbb{R}^d$

$$\begin{aligned} \mathbb{P}(\text{word } i, j \mid \text{topic } h) &= \mathbb{P}(\text{word } i \mid \text{topic } h) \mathbb{P}(\text{word } j \mid \text{topic } h) \\ &= (\mu_h)_i (\mu_h)_j \\ &= (\mu_h \circ \mu_h)_{i,j} \end{aligned}$$

$$\begin{aligned} \mathbb{P}(\text{word } i, j) &= \sum_{h=1}^K \mathbb{P}(\text{topic } h) \mathbb{P}(\text{word } i, j \mid \text{topic } h) \\ &= \sum_{h=1}^K \mu_h (\mu_h \circ \mu_h)_{i,j} \end{aligned}$$

$$\mathbb{P}(\text{word } i_1, i_2, i_3) = \sum_{h=1}^K \mu_h (\mu_h \circ \mu_h \circ \mu_h)_{i_1, i_2, i_3}$$

Note that if x_1, x_2 and x_3 are the one-hot encodings of the first 3 words in a document, then

$$\mathbb{E}[x_1 \circ x_2]_{i_1, i_2} = \mathbb{P}(\text{word } i_1, i_2)$$

$$\mathbb{E}[x_1 \circ x_2 \circ x_3]_{i_1, i_2, i_3} = \mathbb{P}(\text{word } i_1, i_2, i_3)$$

Theorem 3.1 (Anandkumar et al., 2012c) If

$$\begin{aligned} M_2 &:= \mathbb{E}[x_1 \otimes x_2] \\ M_3 &:= \mathbb{E}[x_1 \otimes x_2 \otimes x_3], \end{aligned}$$

then

$$\begin{aligned} M_2 &= \sum_{i=1}^k w_i \mu_i \otimes \mu_i \\ M_3 &= \sum_{i=1}^k w_i \mu_i \otimes \mu_i \otimes \mu_i. \end{aligned} \quad (\text{w}_i = \mu_i)$$

GMM

parameters : $\Theta = \left\{ \underbrace{\mu_1, \dots, \mu_k}_{\substack{\text{centers} \\ \mathbb{R}^d}}, \underbrace{\Sigma_1, \dots, \Sigma_k}_{\substack{\text{Covariance} \\ \mathbb{R}^{d \times d}}}, \underbrace{w_1, \dots, w_k}_{\text{mixing probabilities}} \right\}$

↳ Here we assume each $\Sigma_i = \sigma_i^2 I$ for each $i=1, \dots, k$
 ↳ $\sigma_i \in \mathbb{R}$

Theorem 3.2 (Hsu and Kakade, 2013) Assume $d \geq k$. The variance σ^2 is the smallest eigenvalue of the covariance matrix $\mathbb{E}[x \otimes x] - \mathbb{E}[x] \otimes \mathbb{E}[x]$. Furthermore, if

$$\begin{aligned} M_2 &:= \mathbb{E}[x \otimes x] - \sigma^2 I \\ M_3 &:= \mathbb{E}[x \otimes x \otimes x] - \sigma^2 \sum_{i=1}^d (\mathbb{E}[x] \otimes e_i \otimes e_i + e_i \otimes \mathbb{E}[x] \otimes e_i + e_i \otimes e_i \otimes \mathbb{E}[x]), \end{aligned}$$

then

$$\begin{aligned} M_2 &= \sum_{i=1}^k w_i \mu_i \otimes \mu_i \\ M_3 &= \sum_{i=1}^k w_i \mu_i \otimes \mu_i \otimes \mu_i. \end{aligned}$$

$$\sigma = \sigma_1 = \sigma_2 = \dots = \sigma_k$$

(e_i : i-th vector of the canonical basis)

Summary : In both cases, there exists $M_2 \in \mathbb{R}^{d \times d}$ and a tensor $M_3 \in \mathbb{R}^{d \times d \times d}$
 (which can be estimated from data) satisfying

$$\begin{aligned} M_2 &= \sum_{i=1}^k w_i \mu_i \otimes \mu_i \\ M_3 &= \sum_{i=1}^k w_i \mu_i \otimes \mu_i \otimes \mu_i. \end{aligned} \quad \rightarrow \text{CP decomposition of } M_3 \text{ (under mild conditions, this decomposition unique)}$$

Decomposition of M_2 and M_3 :

$$M_2 = \sum_{i=1}^k w_i \mu_i \otimes \mu_i$$

$$M_3 = \sum_{i=1}^k w_i \mu_i \otimes \mu_i \otimes \mu_i.$$

Observation: . The problem of orthogonal tensor decomposition,

$$T = \sum_{i=1}^k w_i \mu_i \circ \mu_i \circ \mu_i \quad \text{where } \langle \mu_i, \mu_j \rangle = \begin{cases} 1 & \text{if } i=j \\ 0 & \text{o.w.} \end{cases}$$

is easy to solve.

- We can use M_2 to "orthogonalize" M_3 .

