

Guaranteed Learning of Latent Variable Models through Tensor Methods

Furong Huang

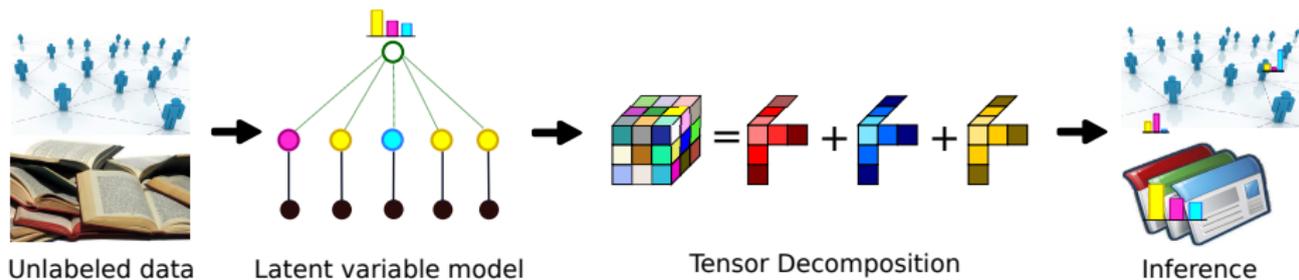
University of Maryland

furongh@cs.umd.edu

ACM SIGMETRICS Tutorial 2018

Tutorial Topic

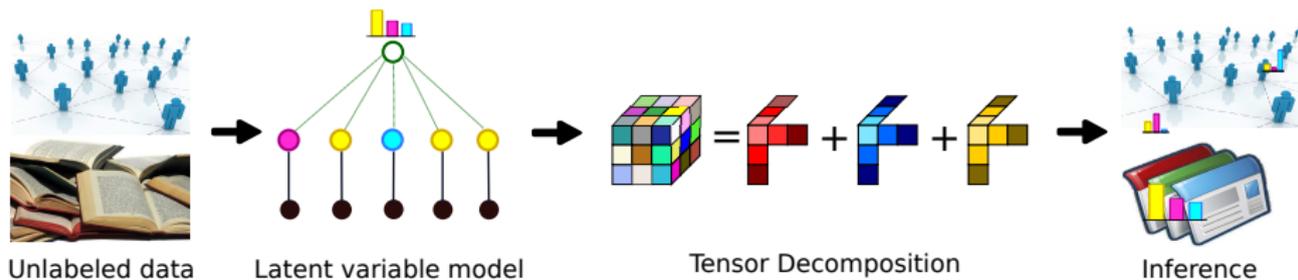
Learning algorithms
for **latent variable models**
based on **decompositions of moment tensors**.



“Method-of-moments” (Pearson, 1894)

Tutorial Topic

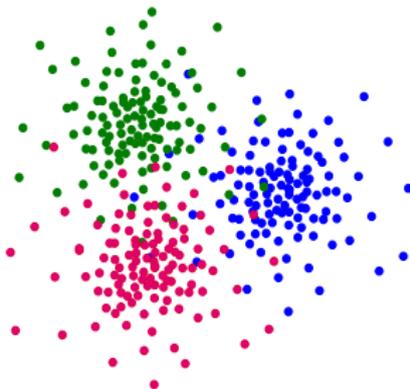
Learning algorithms (parameter estimation)
for **latent variable models**
based on **decompositions of moment tensors**.



“Method-of-moments” (Pearson, 1894)

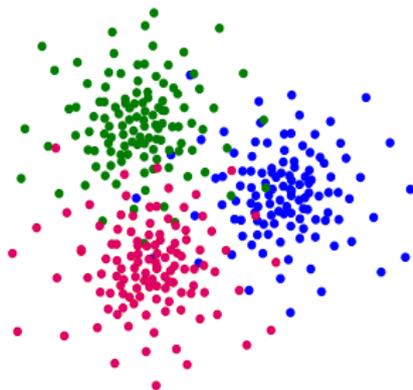
Application 1: Clustering

- Basic operation of grouping data points.
- Hypothesis: each data point belongs to an unknown group.



Application 1: Clustering

- Basic operation of grouping data points.
- Hypothesis: each data point belongs to an unknown group.



Probabilistic/latent variable viewpoint

- The groups represent different distributions. (e.g. Gaussian).
- Each data point is drawn from one of the given distributions. (e.g. Gaussian mixtures).

Application 2: Topic Modeling



Document modeling

- **Observed:** words in document corpus.
- **Hidden:** topics.
- **Goal:** carry out document summarization.

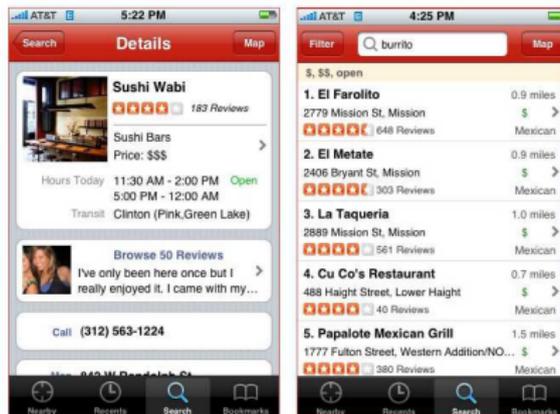
Application 3: Understanding Human Communities



Social Networks

- **Observed:** network of social ties, e.g. friendships, co-authorships
- **Hidden:** groups/communities of social actors.

Application 4: Recommender Systems



Recommender System

- **Observed:** Ratings of users for various products, e.g. yelp reviews.
- **Goal:** Predict new recommendations.
- **Modeling:** Find groups/communities of users and products.

Application 5: Feature Learning

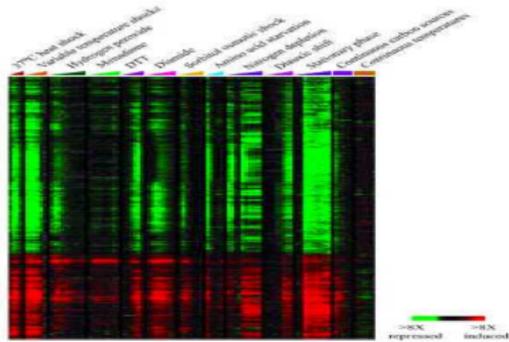


Label	Features				
0	2.1	5.2	0	0	—
1	0	0	2	1	—
1	1.1	0	0	0	—
0	0	0	7	0	—

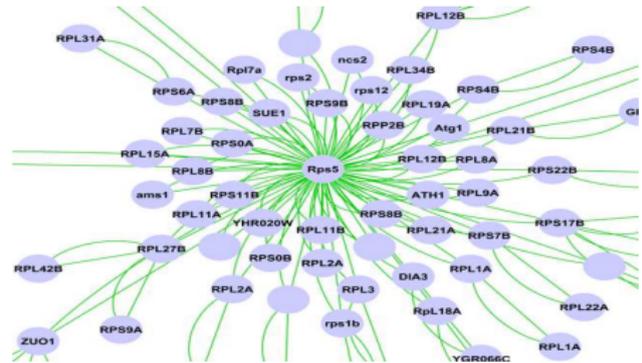
Feature Engineering

- Learn good features/representations for classification tasks, e.g. **image** and **speech recognition**.
- **Sparse** representations, low dimensional hidden structures.

Application 6: Computational Biology



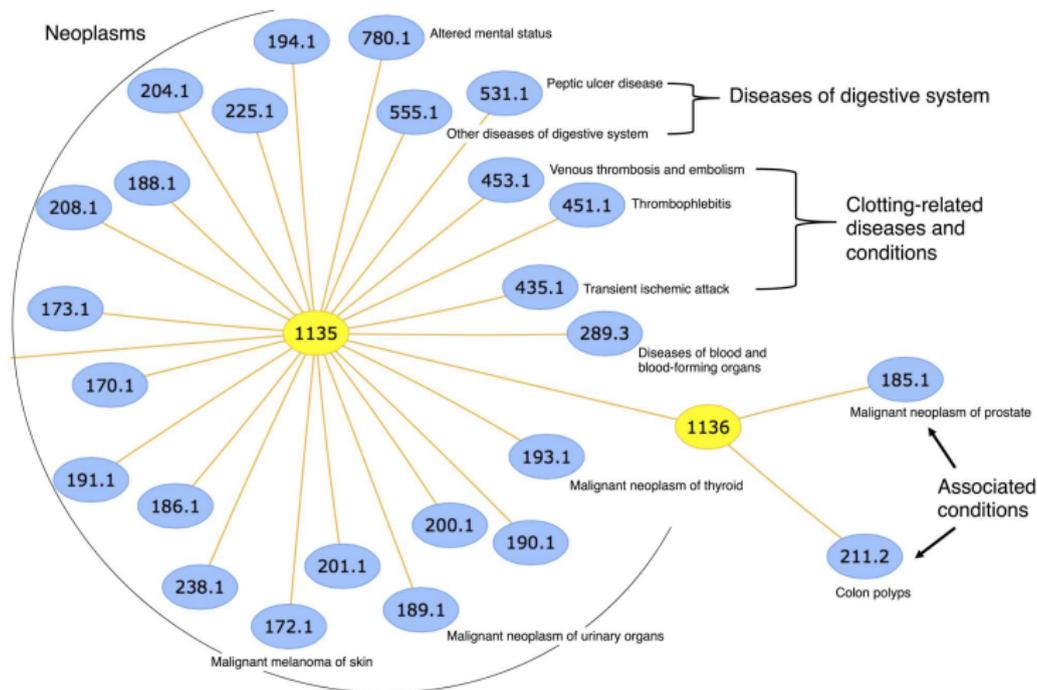
Gasch et al. Mol Biol Cell 2000.



- **Observed:** gene expression levels
- **Goal:** discover gene groups
- **Hidden variables:** regulators controlling gene groups

Application 7: Human Disease Hierarchy Discovery

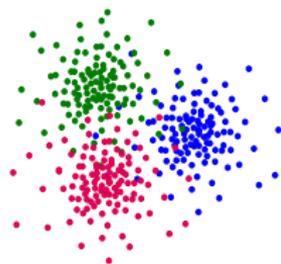
CMS: 1.6 million patients, 168 million diagnostic events, 11 k diseases.



How to model hidden effects?

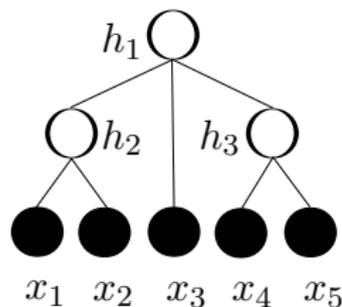
Basic Approach: mixtures/clusters

- Hidden variable h is **categorical**.



Advanced: Probabilistic models

- Hidden variable h has more general distributions.
- Can model mixed memberships.

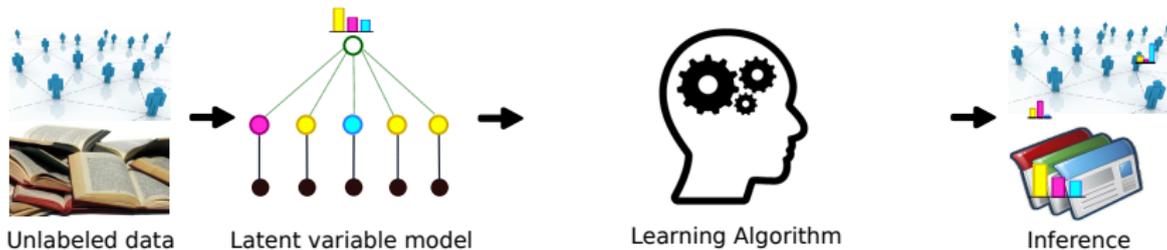


This talk: basic mixture model and some advanced models.

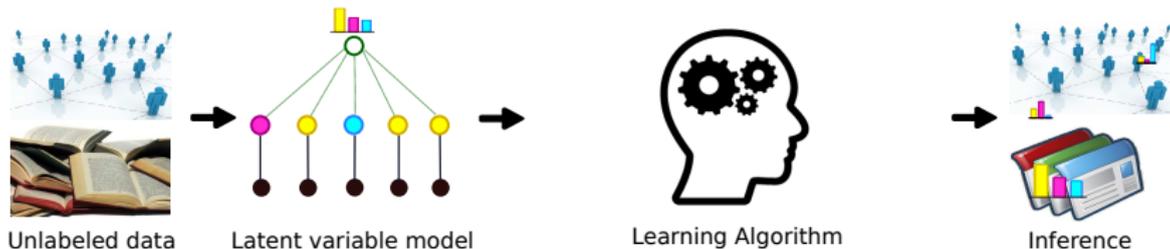
Challenges in Learning

Basic goal in all mentioned applications

Discover hidden structure in data: **unsupervised** learning.



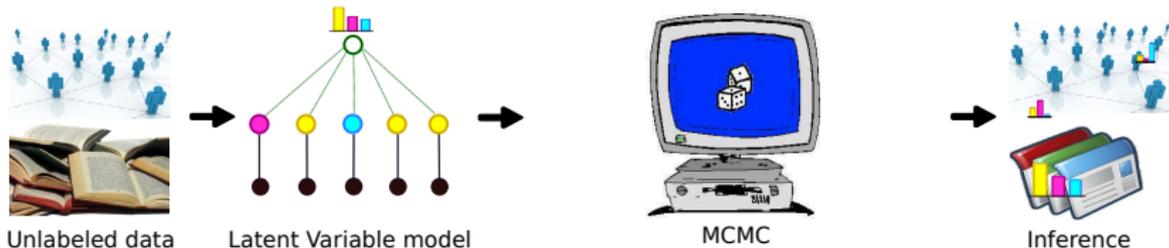
Challenges in Learning – find hidden structure in data



Challenge: Conditions for Identifiability

- Whether can model be identified given **infinite computation and data**?
- Are there **tractable algorithms** under identifiability?

Challenges in Learning – find hidden structure in data



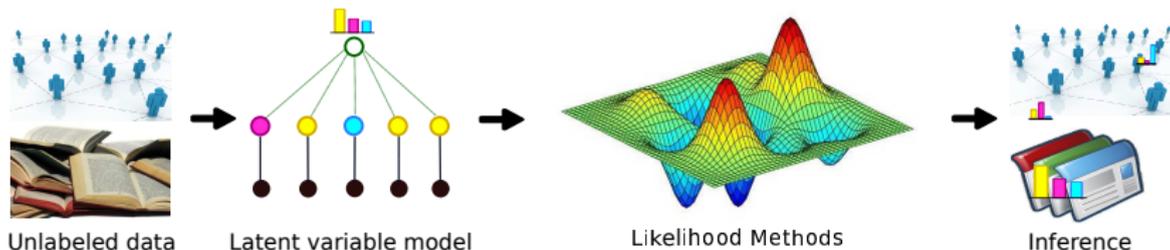
Challenge: Conditions for Identifiability

- Whether can model be identified given **infinite computation and data**?
- Are there **tractable algorithms** under identifiability?

Challenge: Efficient Learning of Latent Variable Models

- MCMC: **random sampling, slow**
Exponential mixing time

Challenges in Learning – find hidden structure in data



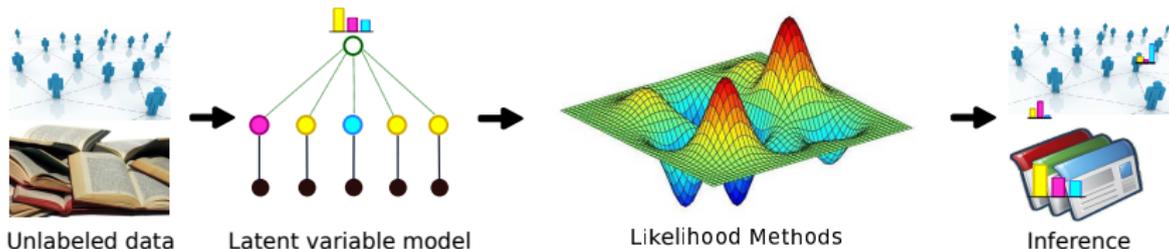
Challenge: Conditions for Identifiability

- Whether can model be identified given **infinite computation and data**?
- Are there **tractable algorithms** under identifiability?

Challenge: Efficient Learning of Latent Variable Models

- MCMC: **random sampling, slow**
Exponential mixing time
- Likelihood: **non-convex, not scalable**
Exponential critical points

Challenges in Learning – find hidden structure in data



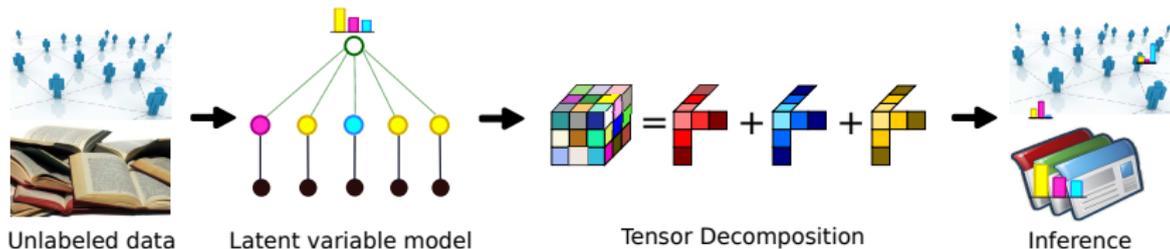
Challenge: Conditions for Identifiability

- Whether can model be identified given **infinite computation and data**?
- Are there **tractable algorithms** under identifiability?

Challenge: Efficient Learning of Latent Variable Models

- MCMC: **random sampling, slow**
Exponential mixing time
- Likelihood: **non-convex, not scalable**
Exponential critical points
- Efficient **computational** and **sample complexities**?

Challenges in Learning – find hidden structure in data



Challenge: Conditions for Identifiability

- Whether can model be identified given **infinite computation and data**?
- Are there **tractable algorithms** under identifiability?

Challenge: Efficient Learning of Latent Variable Models

- MCMC: **random sampling, slow**
Exponential mixing time
- Likelihood: **non-convex, not scalable**
Exponential critical points
- Efficient **computational** and **sample complexities**?

Guaranteed and efficient learning through spectral methods

What this tutorial will cover

Outline

- 1 Introduction
- 2 Motivation: Challenges of MLE for Gaussian Mixtures

What this tutorial will cover

Outline

- 1 Introduction
- 2 Motivation: Challenges of MLE for **Gaussian Mixtures**
- 3 Introduction of **Method of Moments** and **Tensor** Notations

What this tutorial will cover

Outline

- 1 Introduction
- 2 Motivation: Challenges of MLE for **Gaussian Mixtures**
- 3 Introduction of **Method of Moments** and **Tensor** Notations
- 4 Topic Model for Single-topic Documents

What this tutorial will cover

Outline

- 1 Introduction
- 2 Motivation: Challenges of MLE for **Gaussian Mixtures**
- 3 Introduction of **Method of Moments** and **Tensor** Notations
- 4 Topic Model for Single-topic Documents
 - ▶ Identifiability
 - ▶ Parameter recovery via decomposition of exact moments

What this tutorial will cover

Outline

- 1 Introduction
- 2 Motivation: Challenges of MLE for **Gaussian Mixtures**
- 3 Introduction of **Method of Moments** and **Tensor** Notations
- 4 Topic Model for Single-topic Documents
 - ▶ Identifiability
 - ▶ Parameter recovery via decomposition of exact moments
- 5 **Error-tolerant Algorithms for Tensor Decompositions**

What this tutorial will cover

Outline

- 1 Introduction
- 2 Motivation: Challenges of MLE for **Gaussian Mixtures**
- 3 Introduction of **Method of Moments** and **Tensor** Notations
- 4 Topic Model for Single-topic Documents
 - ▶ Identifiability
 - ▶ Parameter recovery via decomposition of exact moments
- 5 **Error-tolerant Algorithms for Tensor Decompositions**
 - ▶ Decomposition for tensors with linearly independent components
 - ▶ Decomposition for tensors with orthogonal components

What this tutorial will cover

Outline

- 1 Introduction
- 2 Motivation: Challenges of MLE for **Gaussian Mixtures**
- 3 Introduction of **Method of Moments** and **Tensor** Notations
- 4 Topic Model for Single-topic Documents
 - ▶ Identifiability
 - ▶ Parameter recovery via decomposition of exact moments
- 5 **Error-tolerant** Algorithms for **Tensor Decompositions**
 - ▶ Decomposition for tensors with linearly independent components
 - ▶ Decomposition for tensors with orthogonal components
- 6 Tensor Decomposition for **Neural Network Compression**
- 7 Conclusion

Outline

- 1 Introduction
- 2 Motivation: Challenges of MLE for Gaussian Mixtures**
- 3 Introduction of Method of Moments and Tensor Notations
- 4 Topic Model for Single-topic Documents
- 5 Algorithms for Tensor Decompositions
- 6 Tensor Decomposition for Neural Network Compression
- 7 Conclusion

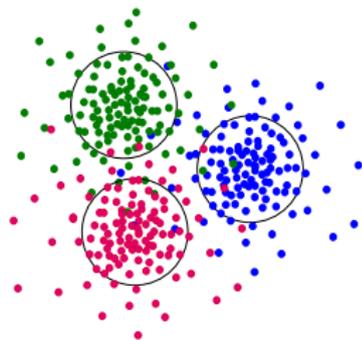
Gaussian Mixture Model

Generative Model

- Samples are comprised of K different Gaussians according to $\text{Cat}(\pi_1, \pi_2, \dots, \pi_K)$
- Each sample is from **one** of the K Gaussians, $\mathcal{N}(\boldsymbol{\mu}_h, \boldsymbol{\Sigma}_h), \forall h \in [K]$

$$H \sim \text{Cat}(\pi_1, \pi_2, \dots, \pi_K)$$

$$\mathbf{X} |_{H=h} \sim \mathcal{N}(\boldsymbol{\mu}_h, \boldsymbol{\Sigma}_h), \quad \forall h \in [K]$$



Gaussian Mixture Model

Generative Model

- Samples are comprised of K different Gaussians according to $\text{Cat}(\pi_1, \pi_2, \dots, \pi_K)$
- Each sample is from **one** of the K Gaussians, $\mathcal{N}(\boldsymbol{\mu}_h, \boldsymbol{\Sigma}_h), \forall h \in [K]$

$$H \sim \text{Cat}(\pi_1, \pi_2, \dots, \pi_K)$$
$$\mathbf{X} |_{H=h} \sim \mathcal{N}(\boldsymbol{\mu}_h, \boldsymbol{\Sigma}_h), \quad \forall h \in [K]$$



Learning Problem

Estimate **mean vector** $\boldsymbol{\mu}_h$, **covariance matrix** $\boldsymbol{\Sigma}_h$, and **mixing weight** $\text{Cat}(\pi_1, \pi_2, \dots, \pi_K)$ of each subpopulation from **unlabeled data**.

Maximum Likelihood Estimator (MLE)

- Data $\{\mathbf{x}_i\}_{i=1}^n$
- Likelihood $\Pr_{\theta}(\text{data}) \stackrel{\text{iid}}{=} \prod_{i=1}^n \Pr_{\theta}(\mathbf{x}_i)$
- Model parameter estimation $\hat{\theta}_{\text{mle}} := \underset{\theta \in \Theta}{\operatorname{argmax}} \log \Pr_{\theta}(\text{data})$
- Latent variable models: some variables are hidden
 - ▶ No “direct” estimators when some variables are hidden
 - ▶ Local optimization via [Expectation-Maximization \(EM\)](#) (Dempster, Laird, & Rubin, 1977)

MLE for Gaussian Mixture Models

Given data $\{\mathbf{x}_i\}_{i=1}^n$ and the number of Gaussian components K , the model parameters to be estimated are $\boldsymbol{\theta} = \{(\boldsymbol{\mu}_h, \boldsymbol{\Sigma}_h, \pi_h)\}_{h=1}^K$.

$\hat{\boldsymbol{\theta}}_{\text{mle}}$ for Gaussian Mixture Models

$$\hat{\boldsymbol{\theta}}_{\text{mle}} := \underset{\boldsymbol{\theta}}{\operatorname{argmax}} \sum_{i=1}^n \log \left(\sum_{h=1}^K \frac{\pi_h}{\det(\boldsymbol{\Sigma}_h)^{1/2}} \exp \left(-\frac{1}{2} (\mathbf{x}_i - \boldsymbol{\mu}_h)^\top \boldsymbol{\Sigma}_h^{-1} (\mathbf{x}_i - \boldsymbol{\mu}_h) \right) \right)$$

- Solving MLE estimator is **NP-hard** (Dasgupta, 2008; Aloise, Deshpande, Hansen, & Popat, 2009; Mahajan, Nimbhorkar, & Varadarajan, 2009; Vattani, 2009; Awasthi, Charikar, Krishnaswamy, & Sinop, 2015).

Consistent Estimator

Definition

Suppose iid samples $\{\mathbf{x}_i\}_{i=1}^n$ are generated by distribution $\Pr_{\theta}(\mathbf{x}_i)$ where the model parameters $\theta \in \Theta$ are unknown. An estimator $\hat{\theta}$ is consistent if

$$\mathbb{E}\|\hat{\theta} - \theta\| \rightarrow 0 \quad \text{as } n \rightarrow \infty$$

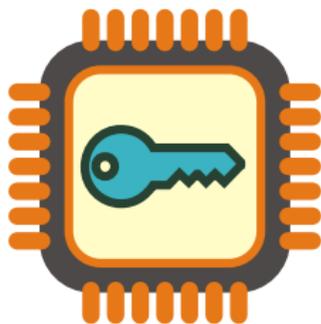
Spherical Gaussian Mixtures $\Sigma_h = I$ (as $n \rightarrow \infty$)

- For $K = 2$ and $\pi_h = 1/2$: **EM is consistent** (Xu, H., & Maleki, 2016; Daskalakis, Tzamos, & Zampetakis, 2016).
- Larger K: **easily trapped in local maxima, far from global max** (Jin, Zhang, Balakrishnan, Wainwright, & Jordan, 2016).
- Practitioners often use EM with many (random) restarts, **but may take a long time to get near the global max.**

Hardness of Parameter Estimation

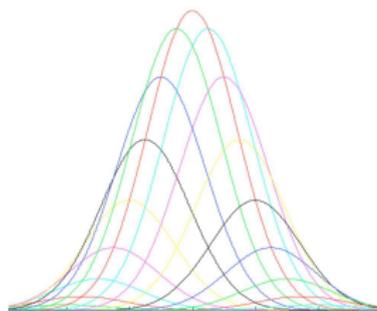
Exponentially difficult computationally or statistically to learn model parameters, even under the parametric setting.

Cryptographic hardness



E.g., Mossel & Roch, 2006

Information-theoretic hardness



E.g., Moitra & Valiant, 2010

May require $2^{\Omega(K)}$ running time or $2^{\Omega(K)}$ sample size.

Ways Around the Hardness

- Separation conditions.

E.g., assume $\min_{i \neq j} \frac{\|\mu_i - \mu_j\|^2}{\sigma_i^2 + \sigma_j^2}$ is sufficiently large.

(Dasgupta, 1999; Arora & Kannan, 2001; Vempala & Wang, 2002; . . .)

- Structural assumptions.

E.g., assume sparsity, separable (anchor words).

(Spielman, Wang & Wright, 2012; Arora, Ge & Moitra, 2012; . . .)

- Non-degeneracy conditions.

E.g., assume μ_1, \dots, μ_K span a K -dimensional space.

This tutorial: statistically and computationally efficient learning algorithms for **non-degenerate** instances via **method-of-moments**.

Outline

- 1 Introduction
- 2 Motivation: Challenges of MLE for Gaussian Mixtures
- 3 Introduction of Method of Moments and Tensor Notations**
- 4 Topic Model for Single-topic Documents
- 5 Algorithms for Tensor Decompositions
- 6 Tensor Decomposition for Neural Network Compression
- 7 Conclusion

Method-of-Moments At A Glance

- 1 Determine **function of model parameters** θ estimatable from observable data:

- ▶ **Moments**

$$\mathbb{E}_{\theta}[f(\mathbf{X})]$$

- 2 Form estimates of moments using data (iid samples $\{\mathbf{x}_i\}_{i=1}^n$):

- ▶ **Empirical Moments**

$$\hat{\mathbb{E}}[f(\mathbf{X})]$$

- 3 Solve the approximate equations for parameters θ :

- ▶ **Moment matching**

$$\mathbb{E}_{\theta}[f(\mathbf{X})] \stackrel{n \rightarrow \infty}{\approx} \hat{\mathbb{E}}[f(\mathbf{X})]$$

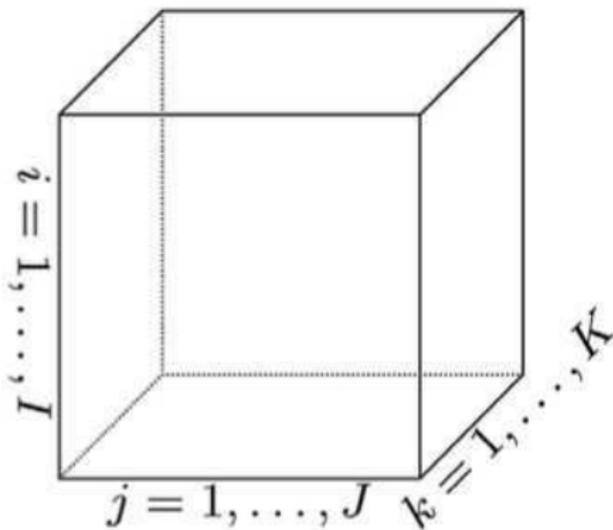
Toy Example

How to estimate Gaussian variable, i.e., $(\boldsymbol{\mu}, \boldsymbol{\Sigma})$,
given iid samples $\{\mathbf{x}_i\}_{i=1}^n \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma}^2)$?

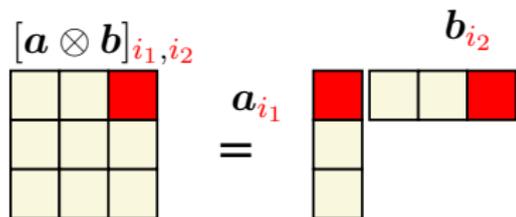
What is a tensor?

Multi-dimensional Array

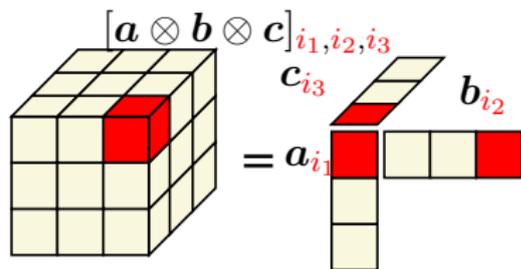
- Tensor - Higher order matrix
- The number of dimensions is called tensor order.



Tensor Product

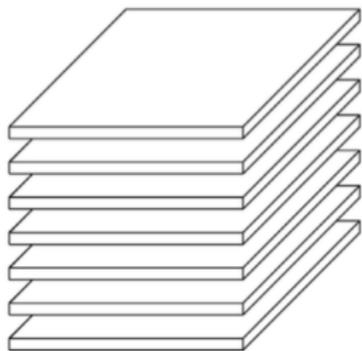

$$[a \otimes b]_{i_1, i_2} = a_{i_1} b_{i_2}$$

- $[a \otimes b]_{i_1, i_2} = a_{i_1} b_{i_2}$
- Rank-1 matrix

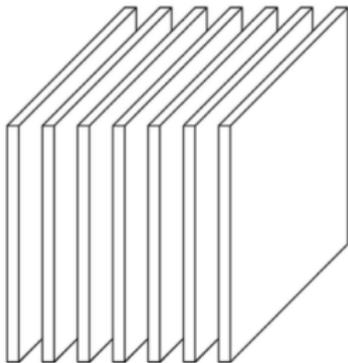

$$[a \otimes b \otimes c]_{i_1, i_2, i_3} = a_{i_1} b_{i_2} c_{i_3}$$

- $[a \otimes b \otimes c]_{i_1, i_2, i_3} = a_{i_1} b_{i_2} c_{i_3}$
- Rank-1 tensor

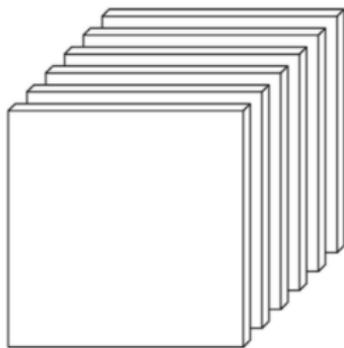
Slices



- Horizontal slices

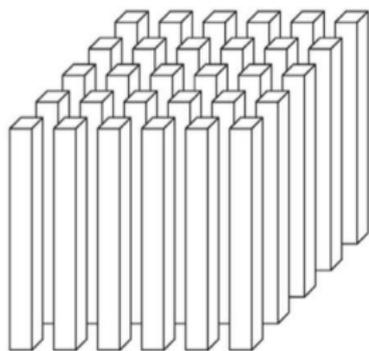


- Lateral slices

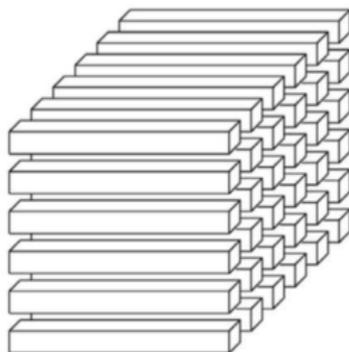


- Frontal slices

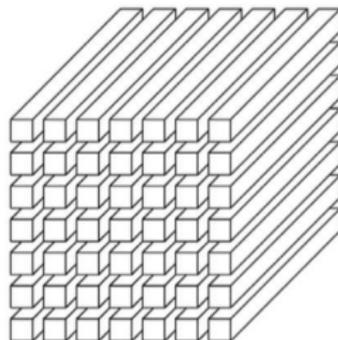
Fiber



- Mode-1 (column) fibers

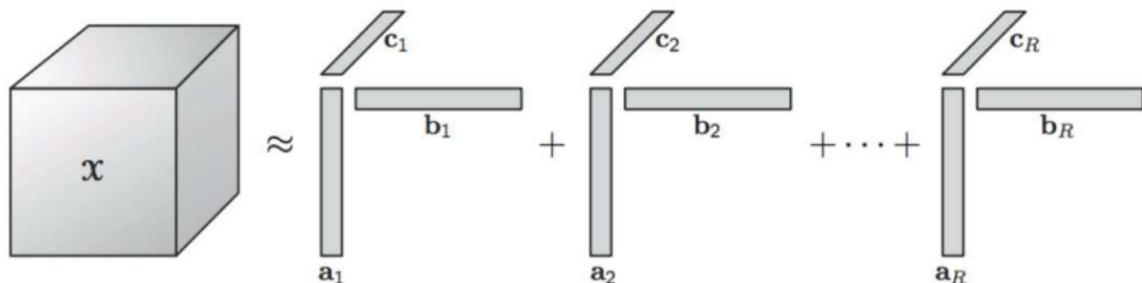


- Mode-2 (row) fibers



- Mode-3 (tube) fibers

CP decomposition



- $\mathcal{X} = \sum_{h=1}^R \mathbf{a}_h \otimes \mathbf{b}_h \otimes \mathbf{c}_h$
- Rank: Minimum number of rank-1 tensors whose sum generates the tensor.

Multi-linear Transform

Multi-linear Operation

If $\mathcal{T} = \sum_{h=1}^R \mathbf{a}_h \otimes \mathbf{b}_h \otimes \mathbf{c}_h$, a multi-linear operation using matrices $(\mathbf{X}, \mathbf{Y}, \mathbf{Z})$ is as follows

$$\mathcal{T}(\mathbf{X}, \mathbf{Y}, \mathbf{Z}) := \sum_{h=1}^K (\mathbf{X}^\top \mathbf{a}_h) \otimes (\mathbf{Y}^\top \mathbf{b}_h) \otimes (\mathbf{Z}^\top \mathbf{c}_h).$$

Similarly for a multi-linear operation using vectors $(\mathbf{x}, \mathbf{y}, \mathbf{z})$

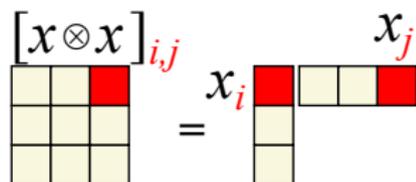
$$\mathcal{T}(\mathbf{x}, \mathbf{y}, \mathbf{z}) := \sum_{h=1}^K (\mathbf{x}^\top \mathbf{a}_h) \otimes (\mathbf{y}^\top \mathbf{b}_h) \otimes (\mathbf{z}^\top \mathbf{c}_h).$$

Tensors in Method of Moments

Matrix: Pair-wise relationship

- Signal or data observed $\mathbf{x} \in \mathbb{R}^d$
- Rank 1 matrix: $[\mathbf{x} \otimes \mathbf{x}]_{i,j} = \mathbf{x}_i \mathbf{x}_j$
- Aggregated pair-wise relationship

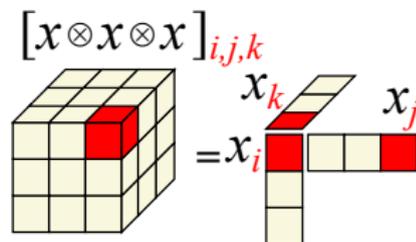
$$\mathcal{M}_2 = \mathbb{E}[\mathbf{x} \otimes \mathbf{x}]$$



Tensor: Triple-wise relationship or higher

- Signal or data observed $\mathbf{x} \in \mathbb{R}^d$
- Rank 1 tensor:
 $[\mathbf{x} \otimes \mathbf{x} \otimes \mathbf{x}]_{i,j,k} = \mathbf{x}_i \mathbf{x}_j \mathbf{x}_k$
- Aggregated triple-wise relationship

$$\mathcal{M}_3 = \mathbb{E}[\mathbf{x} \otimes \mathbf{x} \otimes \mathbf{x}] = \mathbb{E}[\mathbf{x} \otimes^3]$$

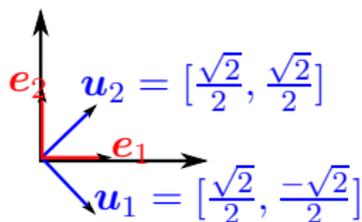


Why are tensors powerful?

Matrix Orthogonal Decomposition

- **Not unique** without eigenvalue gap

$$\begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} = \mathbf{e}_1 \mathbf{e}_1^\top + \mathbf{e}_2 \mathbf{e}_2^\top = \mathbf{u}_1 \mathbf{u}_1^\top + \mathbf{u}_2 \mathbf{u}_2^\top$$



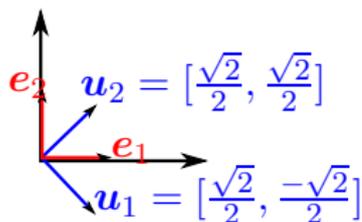
Why are tensors powerful?

Matrix Orthogonal Decomposition

- **Not unique** without eigenvalue gap

$$\begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} = \mathbf{e}_1 \mathbf{e}_1^\top + \mathbf{e}_2 \mathbf{e}_2^\top = \mathbf{u}_1 \mathbf{u}_1^\top + \mathbf{u}_2 \mathbf{u}_2^\top$$

- **Unique** with eigenvalue gap



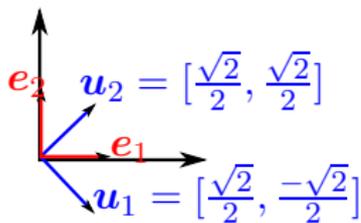
Why are tensors powerful?

Matrix Orthogonal Decomposition

- **Not unique** without eigenvalue gap

$$\begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} = \mathbf{e}_1 \mathbf{e}_1^T + \mathbf{e}_2 \mathbf{e}_2^T = \mathbf{u}_1 \mathbf{u}_1^T + \mathbf{u}_2 \mathbf{u}_2^T$$

- **Unique** with eigenvalue gap



Tensor Orthogonal Decomposition (Harshman, 1970)

- **Unique**: eigenvalue gap not needed

$Tensor = \mathbf{u}_1 \otimes \mathbf{u}_1 \otimes \mathbf{u}_1 + \mathbf{u}_2 \otimes \mathbf{u}_2 \otimes \mathbf{u}_2$

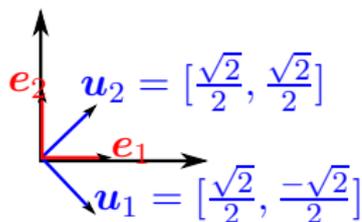
Why are tensors powerful?

Matrix Orthogonal Decomposition

- **Not unique** without eigenvalue gap

$$\begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} = e_1 e_1^T + e_2 e_2^T = u_1 u_1^T + u_2 u_2^T$$

- **Unique** with eigenvalue gap



Tensor Orthogonal Decomposition (Harshman, 1970)

- **Unique**: eigenvalue gap not needed
- Slice of tensor has eigenvalue gap

$$\text{Slice}_i = u_1(i) u_1 \otimes u_1 + u_2(i) u_2 \otimes u_2$$

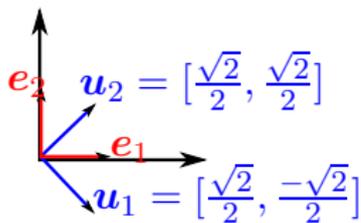
Why are tensors powerful?

Matrix Orthogonal Decomposition

- **Not unique** without eigenvalue gap

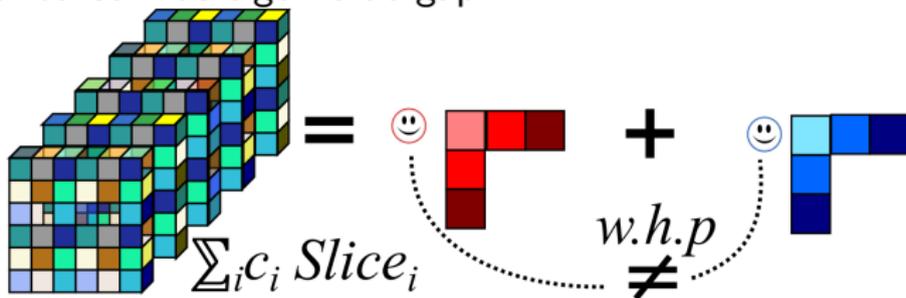
$$\begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} = \mathbf{e}_1 \mathbf{e}_1^\top + \mathbf{e}_2 \mathbf{e}_2^\top = \mathbf{u}_1 \mathbf{u}_1^\top + \mathbf{u}_2 \mathbf{u}_2^\top$$

- **Unique** with eigenvalue gap



Tensor Orthogonal Decomposition (Harshman, 1970)

- **Unique**: eigenvalue gap not needed
- Slice of tensor has eigenvalue gap

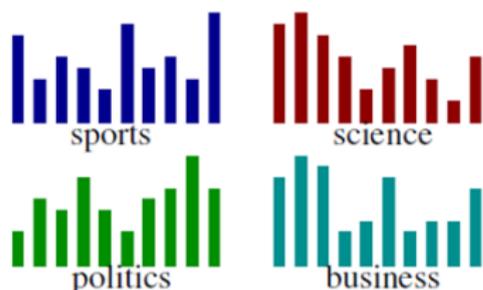


Outline

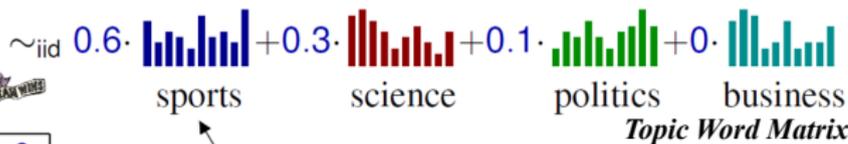
- 1 Introduction
- 2 Motivation: Challenges of MLE for Gaussian Mixtures
- 3 Introduction of Method of Moments and Tensor Notations
- 4 Topic Model for Single-topic Documents**
- 5 Algorithms for Tensor Decompositions
- 6 Tensor Decomposition for Neural Network Compression
- 7 Conclusion

Topic Modeling

General Topic Model (e.g., Latent Dirichlet Allocation)



- K topics
 - ▶ each associated with a **distribution over vocab words** $\{a_h\}_{h=1}^K$
- **Hidden** topic proportion w
 - ▶ per document i , $w^{(i)} \in \Delta^{K-1}$
- Document $\stackrel{iid}{\sim}$ **mixture of topics**



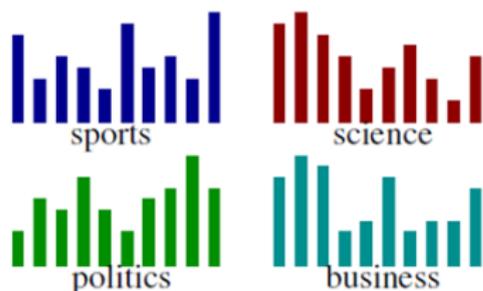
aardvark	0
athlete	3
⋮	⋮
zygote	1

Word Count per Document

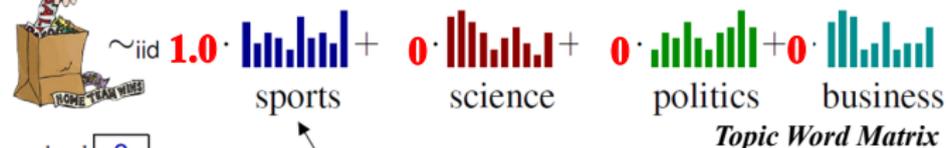


Topic Modeling

Topic Model for Single-topic Documents



- K topics
 - ▶ each associated with a **distribution over vocab words** $\{a_h\}_{h=1}^K$
- **Hidden** topic proportion w
 - ▶ per document i , $w^{(i)} \in \{e_1, \dots, e_K\}$
- Document $\text{iid} \sim a_h$



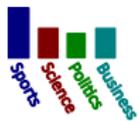
aardvark	0
athlete	3
⋮	⋮
zygote	1

Word Count per Document



Model Parameters of Topic Model for Single-topic Documents

Estimate Topic Proportion



- Topic proportion $\mathbf{w} = [w_1, \dots, w_K]$

$$w_h = \mathbb{P}[\text{topic of word} = h]$$

Estimate Topic Word Matrix



- Topic-word matrix $\mathbf{A} = [\mathbf{a}_1, \dots, \mathbf{a}_K]$

$$A_{jh} = \mathbb{P}[\text{word} = e_j | \text{topic} = h]$$

- **Goal:** to estimate model parameters $\{(\mathbf{a}_h, w_h)\}_{h=1}^K$, given iid samples of n documents (word count $\{\mathbf{c}^{(i)}\}_{i=1}^n$)
- **Frequency vector** $\mathbf{x}^{(i)} = \frac{\mathbf{c}^{(i)}}{L}$, the length of document is $L = \sum_j \mathbf{c}_j^{(i)}$

Moment Matching

Nondegenerate model (linearly independent topic-word matrix)

	Sports	Science	Politics	Business
play	light blue	light green	light yellow	light orange
game	medium blue	medium green	medium yellow	medium orange
season	dark blue	dark green	dark yellow	dark orange

- Generative process:

Choose $h \sim \text{Cat}(w_1, \dots, w_K)$

Generate L words $\sim \mathbf{a}_h$

- $\mathbb{E}[\mathbf{x}] = \sum_{h=1}^K \mathbb{P}[\text{topic} = h] \mathbb{E}[\mathbf{x} | \text{topic} = h]$

- $\mathbb{E}[\mathbf{x} | \text{topic} = h] = \sum_j \mathbb{P}[\text{word} = e_j | \text{topic} = h] e_j = \mathbf{a}_h$

Moment Matching

Nondegenerate model (linearly independent topic-word matrix)



	Sports	Science	Politics	Business
play	light blue	light green	light yellow	light orange
game	medium blue	medium green	medium yellow	medium orange
season	dark blue	dark green	dark yellow	dark orange

- Generative process:

Choose $h \sim \text{Cat}(w_1, \dots, w_K)$

Generate L words $\sim \mathbf{a}_h$

- $\mathbb{E}[\mathbf{x}] = \sum_{h=1}^K \mathbb{P}[\text{topic} = h] \mathbb{E}[\mathbf{x} | \text{topic} = h] = \sum_{h=1}^K w_h \mathbf{a}_h$
- $\mathbb{E}[\mathbf{x} | \text{topic} = h] = \sum_j \mathbb{P}[\text{word} = e_j | \text{topic} = h] e_j = \mathbf{a}_h$

Identifiability: how long must the documents be?

Nondegenerate model (linearly independent topic-word matrix)



- Generative process:

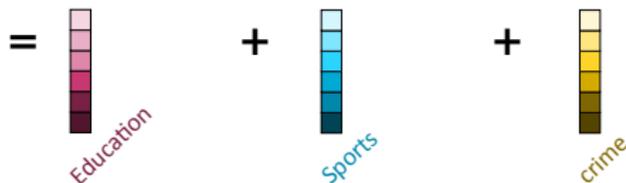
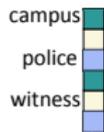
Choose $h \sim \text{Cat}(w_1, \dots, w_K)$

Generate L words $\sim \mathbf{a}_h$

- $\mathbb{E}[\mathbf{x}] = \sum_{h=1}^K \mathbb{P}[\text{topic} = h] \mathbb{E}[\mathbf{x} | \text{topic} = h] = \sum_{h=1}^K w_h \mathbf{a}_h$
- $\mathbb{E}[\mathbf{x} | \text{topic} = h] = \sum_j \mathbb{P}[\text{word} = e_j | \text{topic} = h] e_j = \mathbf{a}_h$

M_1 : Distribution of words (\widehat{M}_1 : Occurrence frequency of words)

$$M_1 = \mathbb{E}[\mathbf{x}] = \sum_h w_h \mathbf{a}_h; \quad \widehat{M}_1 = \frac{1}{n} \sum_{i=1}^n \mathbf{x}^{(i)}$$



Identifiability: how long must the documents be?

Nondegenerate model (linearly independent topic-word matrix)



- Generative process:

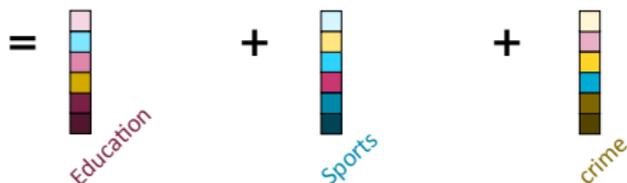
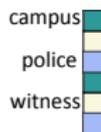
Choose $h \sim \text{Cat}(w_1, \dots, w_K)$

Generate L words $\sim \mathbf{a}_h$

- $\mathbb{E}[\mathbf{x}] = \sum_{h=1}^K \mathbb{P}[\text{topic} = h] \mathbb{E}[\mathbf{x} | \text{topic} = h] = \sum_{h=1}^K w_h \mathbf{a}_h$
- $\mathbb{E}[\mathbf{x} | \text{topic} = h] = \sum_j \mathbb{P}[\text{word} = e_j | \text{topic} = h] e_j = \mathbf{a}_h$

M_1 : Distribution of words (\widehat{M}_1 : Occurrence frequency of words)

$$M_1 = \mathbb{E}[\mathbf{x}] = \sum_h w_h \mathbf{a}_h; \quad \widehat{M}_1 = \frac{1}{n} \sum_{i=1}^n \mathbf{x}^{(i)}$$



No unique decomposition of vectors

Identifiability: how long must the documents be?

Nondegenerate model (linearly independent topic-word matrix)



- Generative process:

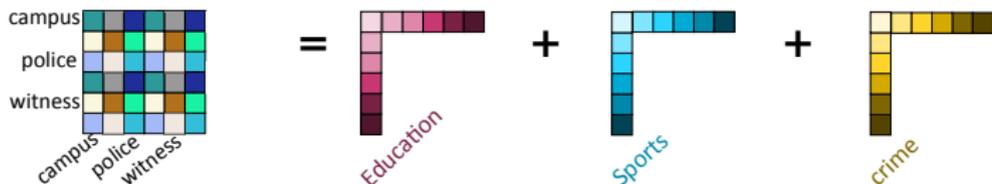
Choose $h \sim \text{Cat}(w_1, \dots, w_K)$

Generate L words $\sim \mathbf{a}_h$

- $\mathbb{E}[\mathbf{x}] = \sum_{h=1}^K \mathbb{P}[\text{topic} = h] \mathbb{E}[\mathbf{x} | \text{topic} = h] = \sum_{h=1}^K w_h \mathbf{a}_h$
- $\mathbb{E}[\mathbf{x} | \text{topic} = h] = \sum_j \mathbb{P}[\text{word} = e_j | \text{topic} = h] e_j = \mathbf{a}_h$

M_2 : Distribution of word pairs (\widehat{M}_2 : Co-occurrence of word pairs)

$$M_2 = \mathbb{E}[\mathbf{x} \otimes \mathbf{x}] = \sum_h w_h \mathbf{a}_h \otimes \mathbf{a}_h; \quad \widehat{M}_2 = \frac{1}{n} \sum_{i=1}^n \mathbf{x}^{(i)} \otimes \mathbf{x}^{(i)}$$



Identifiability: how long must the documents be?

Nondegenerate model (linearly independent topic-word matrix)



- Generative process:

Choose $h \sim \text{Cat}(w_1, \dots, w_K)$

Generate L words $\sim \mathbf{a}_h$

- $\mathbb{E}[\mathbf{x}] = \sum_{h=1}^K \mathbb{P}[\text{topic} = h] \mathbb{E}[\mathbf{x} | \text{topic} = h] = \sum_{h=1}^K w_h \mathbf{a}_h$
- $\mathbb{E}[\mathbf{x} | \text{topic} = h] = \sum_j \mathbb{P}[\text{word} = e_j | \text{topic} = h] e_j = \mathbf{a}_h$

M_2 : Distribution of word pairs (\widehat{M}_2 : Co-occurrence of word pairs)

$$M_2 = \mathbb{E}[\mathbf{x} \otimes \mathbf{x}] = \sum_h w_h \mathbf{a}_h \otimes \mathbf{a}_h; \quad \widehat{M}_2 = \frac{1}{n} \sum_{i=1}^n \mathbf{x}^{(i)} \otimes \mathbf{x}^{(i)}$$

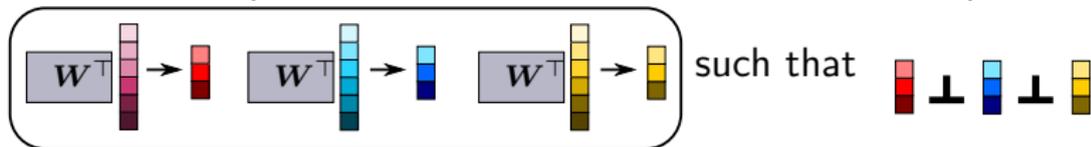


Matrix decomposition recovers subspace, not actual model

Identifiability: how long must the documents be?

Nondegenerate model (linearly independent topic-word matrix)

Find a W



M_2 : Distribution of word pairs (\widehat{M}_2 : Co-occurrence of word pairs)

$$M_2 = \mathbb{E}[\mathbf{x} \otimes \mathbf{x}] = \sum_h w_h \mathbf{a}_h \otimes \mathbf{a}_h; \quad \widehat{M}_2 = \frac{1}{n} \sum_{i=1}^n \mathbf{x}^{(i)} \otimes \mathbf{x}^{(i)}$$

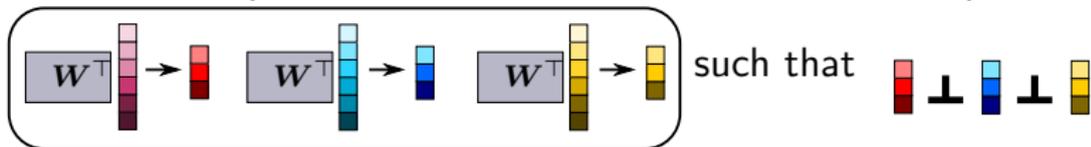


Many such W 's, find one such that $\mathbf{v}_h = W^T \mathbf{a}_h$ orthogonal

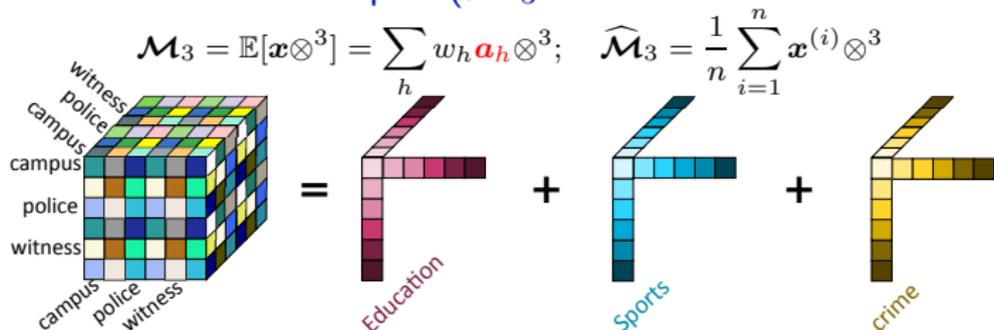
Identifiability: how long must the documents be?

Nondegenerate model (linearly independent topic-word matrix)

Know a W



\mathcal{M}_3 : Distribution of word triples ($\widehat{\mathcal{M}}_3$: Co-occurrence of word triples)

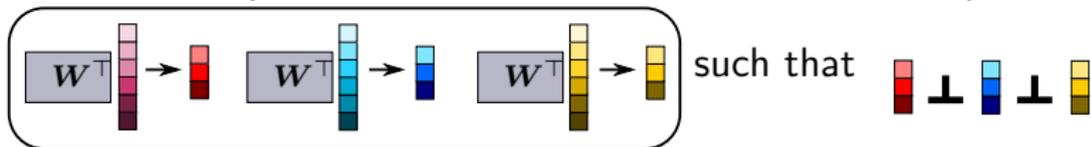


Orthogonalize the tensor, project data with W : $\mathcal{M}_3(W, W, W)$

Identifiability: how long must the documents be?

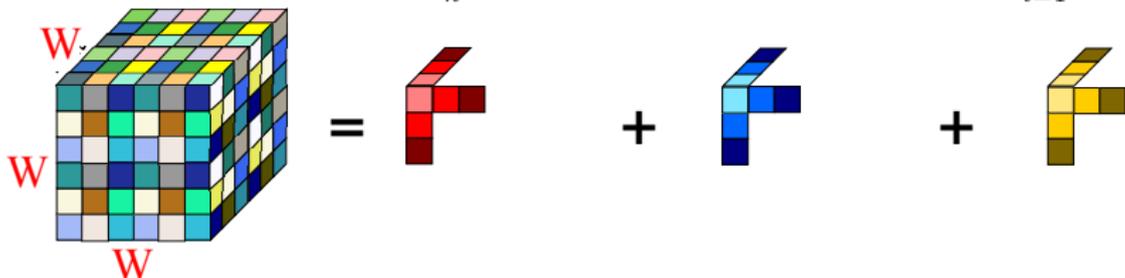
Nondegenerate model (linearly independent topic-word matrix)

Know a W



\mathcal{M}_3 : Distribution of word triples ($\widehat{\mathcal{M}}_3$: Co-occurrence of word triples)

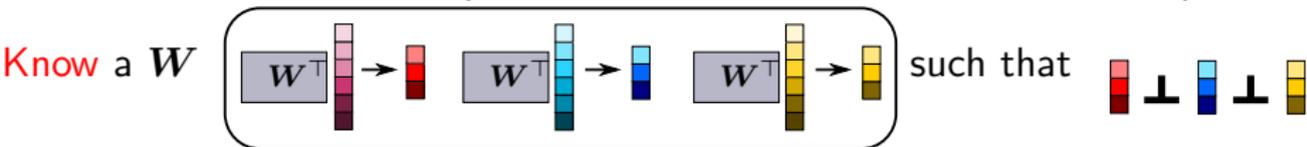
$$\mathcal{M}_3(W, W, W) = \mathbb{E}[(W^\top x) \otimes^3] = \sum_h w_h (W^\top a_h) \otimes^3; \quad \widehat{\mathcal{M}}_3(W, W, W) = \frac{1}{n} \sum_{i=1}^n (W^\top x^{(i)}) \otimes^3$$



Unique orthogonal tensor decomposition $\{\widehat{v}_h\}_{h=1}^K$

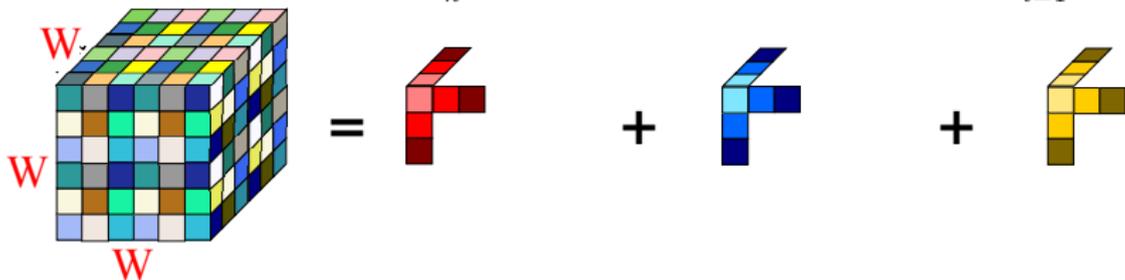
Identifiability: how long must the documents be?

Nondegenerate model (linearly independent topic-word matrix)



\mathcal{M}_3 : Distribution of word triples ($\widehat{\mathcal{M}}_3$: Co-occurrence of word triples)

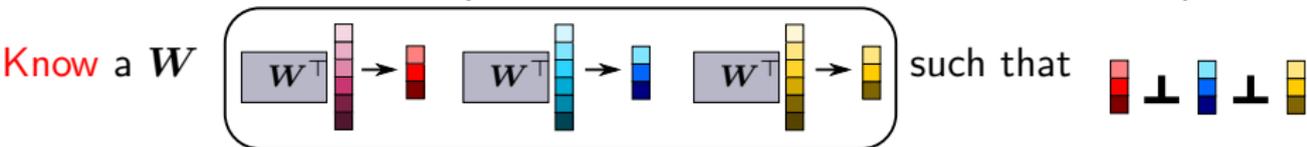
$$\mathcal{M}_3(W, W, W) = \mathbb{E}[(W^\top x) \otimes^3] = \sum_h w_h (W^\top a_h) \otimes^3; \quad \widehat{\mathcal{M}}_3(W, W, W) = \frac{1}{n} \sum_{i=1}^n (W^\top x^{(i)}) \otimes^3$$



$$\text{Model parameter estimation: } \widehat{a}_h = (W^\top)^\dagger \widehat{v}_h$$

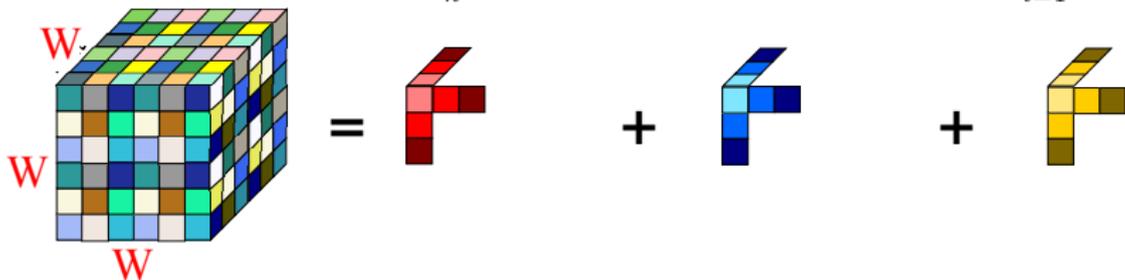
Identifiability: how long must the documents be?

Nondegenerate model (linearly independent topic-word matrix)



\mathcal{M}_3 : Distribution of word triples ($\widehat{\mathcal{M}}_3$: Co-occurrence of word triples)

$$\mathcal{M}_3(W, W, W) = \mathbb{E}[(W^\top x) \otimes^3] = \sum_h w_h (W^\top a_h) \otimes^3; \quad \widehat{\mathcal{M}}_3(W, W, W) = \frac{1}{n} \sum_{i=1}^n (W^\top x^{(i)}) \otimes^3$$



$L \geq 3$: Learning Topic Models through Matrix/Tensor Decomposition

Take Away Message

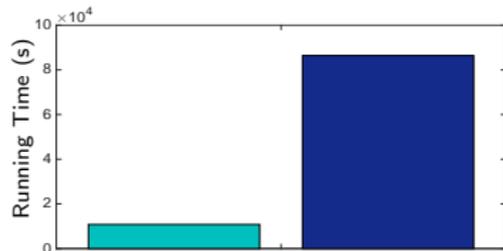
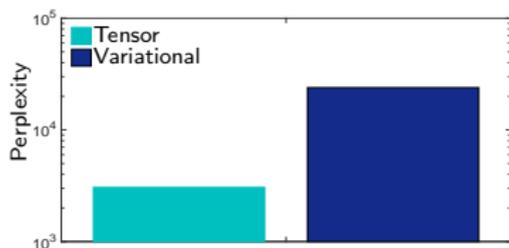
- Consider topic models satisfying **linear independent word distributions** under different **topics**.
- Parameters of topic model for single-topic documents can be efficiently recovered from **distribution of three-word documents**.
 - ▶ Distribution of three-word documents (word triples)

$$M_3 = \mathbb{E}[\mathbf{x} \otimes \mathbf{x} \otimes \mathbf{x}] = \sum_h w_h \mathbf{a}_h \otimes \mathbf{a}_h \otimes \mathbf{a}_h$$

- ▶ \widehat{M}_3 : Co-occurrence of word triples
- Two-word documents are not sufficient for identifiability.

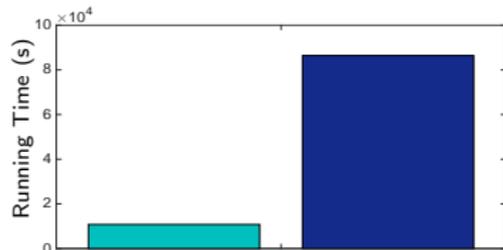
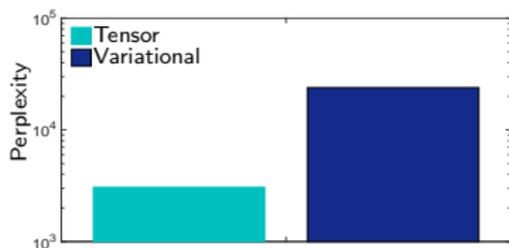
Tensor Methods Compared with Variational Inference

Learning Topics from PubMed on Spark: 8 million docs



Tensor Methods Compared with Variational Inference

Learning Topics from PubMed on Spark: 8 million docs

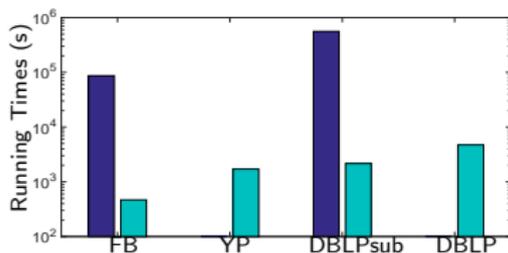
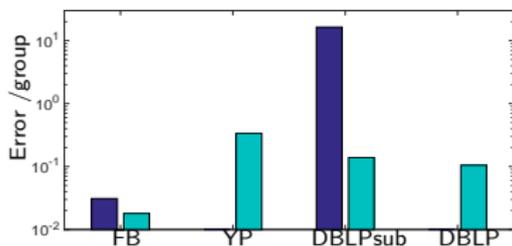


Learning Communities from Graph Connectivity

Facebook: $n \sim 20k$ Yelp: $n \sim 40k$

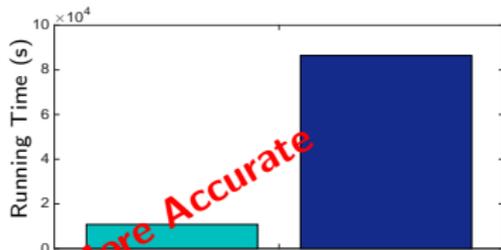
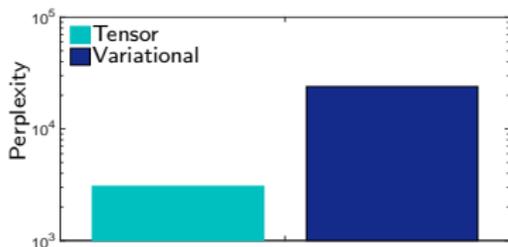
DBLPsub: $n \sim 0.1m$

DBLP: $n \sim 1m$



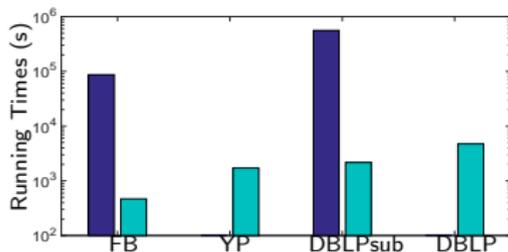
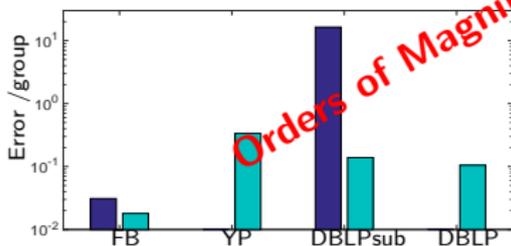
Tensor Methods Compared with Variational Inference

Learning Topics from PubMed on Spark: 8 million docs



Learning Communities from Graph Connectivity

Facebook: $n \sim 20k$ Yelp: $n \sim 40k$ DBLPsub: $n \sim 0.1m$ DBLP: $n \sim 1m$



"Online Tensor Methods for Learning Latent Variable Models", F. Huang, U. Niranjan, M. Hakeem, A. Anandkumar, JMLR14.

"Tensor Methods on Apache Spark", by F. Huang, A. Anandkumar, Oct. 2015.

Outline

- 1 Introduction
- 2 Motivation: Challenges of MLE for Gaussian Mixtures
- 3 Introduction of Method of Moments and Tensor Notations
- 4 Topic Model for Single-topic Documents
- 5 Algorithms for Tensor Decompositions**
- 6 Tensor Decomposition for Neural Network Compression
- 7 Conclusion

Jennrich's Algorithm (Simplified)

Task: Given tensor $\mathcal{T} = \sum_{h=1}^K \mu_h \otimes^3$ with linearly independent components $\{\mu_h\}_{h=1}^K$, find the components (up to scaling).

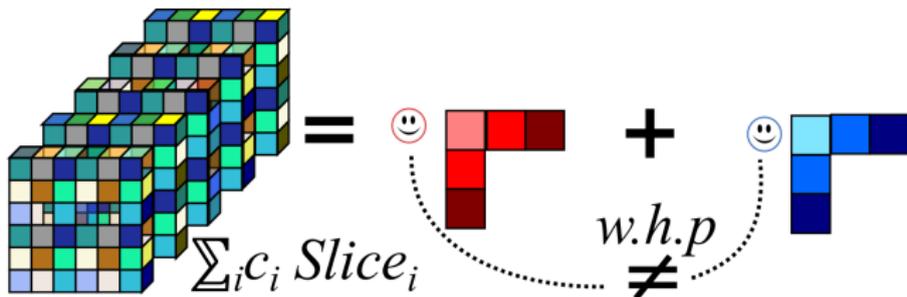
$Tensor = u_1 \otimes u_1 \otimes u_1 + u_2 \otimes u_2 \otimes u_2$

Jennrich's Algorithm (Simplified)

Task: Given tensor $\mathcal{T} = \sum_{h=1}^K \mu_h \otimes^3$ with linearly independent components $\{\mu_h\}_{h=1}^K$, find the components (up to scaling).

Properties of Tensor Slices

- Linear combination of slices $\mathcal{T}(I, I, c) = \sum_h \langle \mu_h, c \rangle \mu_h \otimes \mu_h$

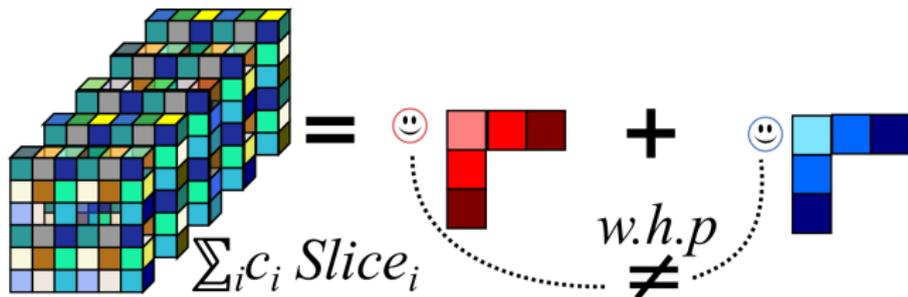


Jennrich's Algorithm (Simplified)

Task: Given tensor $\mathcal{T} = \sum_{h=1}^K \mu_h \otimes^3$ with linearly independent components $\{\mu_h\}_{h=1}^K$, find the components (up to scaling).

Properties of Tensor Slices

- Linear combination of slices $\mathcal{T}(I, I, c) = \sum_h \langle \mu_h, c \rangle \mu_h \otimes \mu_h$



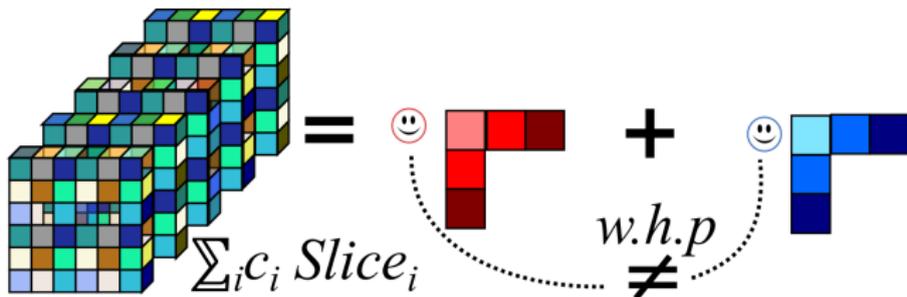
Intuitions for Jennrich's Algorithm

Jennrich's Algorithm (Simplified)

Task: Given tensor $\mathcal{T} = \sum_{h=1}^K \mu_h \otimes^3$ with linearly independent components $\{\mu_h\}_{h=1}^K$, find the components (up to scaling).

Properties of Tensor Slices

- Linear combination of slices $\mathcal{T}(I, I, c) = \sum_h \langle \mu_h, c \rangle \mu_h \otimes \mu_h$



Intuitions for Jennrich's Algorithm

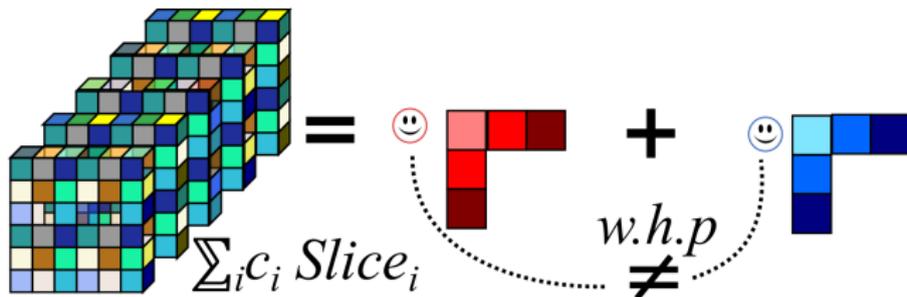
Linear comb. of slices of a tensor share the same set of eigenvectors

Jennrich's Algorithm (Simplified)

Task: Given tensor $\mathcal{T} = \sum_{h=1}^K \mu_h \otimes^3$ with linearly independent components $\{\mu_h\}_{h=1}^K$, find the components (up to scaling).

Properties of Tensor Slices

- Linear combination of slices $\mathcal{T}(I, I, c) = \sum_h \langle \mu_h, c \rangle \mu_h \otimes \mu_h$



Intuitions for Jennrich's Algorithm

Linear comb. of slices of a tensor share the same set of eigenvectors

The shared eigenvectors are tensor components $\{\mu_h\}_{h=1}^K$

Jennrich's Algorithm (Simplified)

Task: Given tensor $\mathcal{T} = \sum_{h=1}^K \mu_h \otimes^3$ with linearly independent components $\{\mu_h\}_{h=1}^K$, find the components (up to scaling).

Algorithm Jennrich's Algorithm

Require: Tensor $\mathcal{T} \in \mathbb{R}^{d \times d \times d}$

Ensure: Components $\{\hat{\mu}_h\}_{h=1}^K \stackrel{\text{a.s.}}{=} \{\mu_h\}_{h=1}^K$

- 1: Sample \mathbf{c} and \mathbf{c}' **independently & uniformly at random** from S^{d-1}
 - 2: Return $\{\hat{\mu}_h\}_{h=1}^K \leftarrow$ eigenvectors of $(\mathcal{T}(\mathbf{I}, \mathbf{I}, \mathbf{c})\mathcal{T}(\mathbf{I}, \mathbf{I}, \mathbf{c}')^\dagger)$
-

Jennrich's Algorithm (Simplified)

Task: Given tensor $\mathcal{T} = \sum_{h=1}^K \boldsymbol{\mu}_h \otimes^3$ with linearly independent components $\{\boldsymbol{\mu}_h\}_{h=1}^K$, find the components (up to scaling).

Algorithm Jennrich's Algorithm

Require: Tensor $\mathcal{T} \in \mathbb{R}^{d \times d \times d}$

Ensure: Components $\{\hat{\boldsymbol{\mu}}_h\}_{h=1}^K \stackrel{\text{a.s.}}{\equiv} \{\boldsymbol{\mu}_h\}_{h=1}^K$

- 1: Sample \mathbf{c} and \mathbf{c}' **independently & uniformly at random** from S^{d-1}
 - 2: Return $\{\hat{\boldsymbol{\mu}}_h\}_{h=1}^K \leftarrow$ eigenvectors of $(\mathcal{T}(\mathbf{I}, \mathbf{I}, \mathbf{c})\mathcal{T}(\mathbf{I}, \mathbf{I}, \mathbf{c}')^\dagger)$
-

Consistency of Jennrich's Algorithm?

Estimators $\{\hat{\boldsymbol{\mu}}_h\}_{h=1}^K \equiv$ unknown components $\{\boldsymbol{\mu}_h\}_{h=1}^K$ (up to scaling)?

Analysis of Consistency of Jennrich's algorithm

Recall: Linear comb. of slices share eigenvectors $\{\boldsymbol{\mu}_h\}_{h=1}^K$,

i.e.,

$$\mathcal{T}(I, I, \mathbf{c})\mathcal{T}(I, I, \mathbf{c}')^\dagger \stackrel{\text{a.s.}}{=} \mathbf{U}\mathbf{D}_c\mathbf{U}^\top (\mathbf{U}^\top)^\dagger \mathbf{D}_{c'}^{-1}\mathbf{U}^\dagger \stackrel{\text{a.s.}}{=} \mathbf{U}(\mathbf{D}_c\mathbf{D}_{c'}^{-1})\mathbf{U}^\dagger,$$

where $\mathbf{U} = [\boldsymbol{\mu}_1 | \dots | \boldsymbol{\mu}_K]$ are the linearly independent tensor components and $\mathbf{D}_c = \text{Diag}(\langle \boldsymbol{\mu}_1, \mathbf{c} \rangle, \dots, \langle \boldsymbol{\mu}_K, \mathbf{c} \rangle)$ is diagonal.

Analysis of Consistency of Jennrich's algorithm

Recall: Linear comb. of slices share eigenvectors $\{\boldsymbol{\mu}_h\}_{h=1}^K$,

i.e.,

$$\mathcal{T}(I, I, \mathbf{c})\mathcal{T}(I, I, \mathbf{c}')^\dagger \stackrel{\text{a.s.}}{=} \mathbf{U}\mathbf{D}_c\mathbf{U}^\top (\mathbf{U}^\top)^\dagger \mathbf{D}_{c'}^{-1}\mathbf{U}^\dagger \stackrel{\text{a.s.}}{=} \mathbf{U}(\mathbf{D}_c\mathbf{D}_{c'}^{-1})\mathbf{U}^\dagger,$$

where $\mathbf{U} = [\boldsymbol{\mu}_1 | \dots | \boldsymbol{\mu}_K]$ are the linearly independent tensor components and $\mathbf{D}_c = \text{Diag}(\langle \boldsymbol{\mu}_1, \mathbf{c} \rangle, \dots, \langle \boldsymbol{\mu}_K, \mathbf{c} \rangle)$ is diagonal.

By linear independence of $\{\boldsymbol{\mu}_i\}_{i=1}^K$ and random choice of \mathbf{c} and \mathbf{c}' :

- 1 \mathbf{U} has rank K ;

Analysis of Consistency of Jennrich's algorithm

Recall: Linear comb. of slices share eigenvectors $\{\boldsymbol{\mu}_h\}_{h=1}^K$,

i.e.,

$$\mathcal{T}(I, I, \mathbf{c})\mathcal{T}(I, I, \mathbf{c}')^\dagger \stackrel{\text{a.s.}}{=} \mathbf{U}\mathbf{D}_c\mathbf{U}^\top (\mathbf{U}^\top)^\dagger \mathbf{D}_{c'}^{-1}\mathbf{U}^\dagger \stackrel{\text{a.s.}}{=} \mathbf{U}(\mathbf{D}_c\mathbf{D}_{c'}^{-1})\mathbf{U}^\dagger,$$

where $\mathbf{U} = [\boldsymbol{\mu}_1 | \dots | \boldsymbol{\mu}_K]$ are the linearly independent tensor components and $\mathbf{D}_c = \text{Diag}(\langle \boldsymbol{\mu}_1, \mathbf{c} \rangle, \dots, \langle \boldsymbol{\mu}_K, \mathbf{c} \rangle)$ is diagonal.

By linear independence of $\{\boldsymbol{\mu}_i\}_{i=1}^K$ and random choice of \mathbf{c} and \mathbf{c}' :

- 1 \mathbf{U} has rank K ;
- 2 \mathbf{D}_c and $\mathbf{D}_{c'}$ are invertible (a.s.);

Analysis of Consistency of Jennrich's algorithm

Recall: Linear comb. of slices share eigenvectors $\{\boldsymbol{\mu}_h\}_{h=1}^K$,

i.e.,

$$\mathcal{T}(I, I, \mathbf{c})\mathcal{T}(I, I, \mathbf{c}')^\dagger \stackrel{\text{a.s.}}{=} \mathbf{U}\mathbf{D}_c\mathbf{U}^\top (\mathbf{U}^\top)^\dagger \mathbf{D}_{c'}^{-1}\mathbf{U}^\dagger \stackrel{\text{a.s.}}{=} \mathbf{U}(\mathbf{D}_c\mathbf{D}_{c'}^{-1})\mathbf{U}^\dagger,$$

where $\mathbf{U} = [\boldsymbol{\mu}_1 | \dots | \boldsymbol{\mu}_K]$ are the linearly independent tensor components and $\mathbf{D}_c = \text{Diag}(\langle \boldsymbol{\mu}_1, \mathbf{c} \rangle, \dots, \langle \boldsymbol{\mu}_K, \mathbf{c} \rangle)$ is diagonal.

By linear independence of $\{\boldsymbol{\mu}_i\}_{i=1}^K$ and random choice of \mathbf{c} and \mathbf{c}' :

- 1 \mathbf{U} has rank K ;
- 2 \mathbf{D}_c and $\mathbf{D}_{c'}$ are invertible (a.s.);
- 3 Diagonal entries of $\mathbf{D}_c\mathbf{D}_{c'}^{-1}$ are distinct (a.s.);

Analysis of Consistency of Jennrich's algorithm

Recall: Linear comb. of slices share eigenvectors $\{\boldsymbol{\mu}_h\}_{h=1}^K$,

i.e.,

$$\mathcal{T}(\mathbf{I}, \mathbf{I}, \mathbf{c})\mathcal{T}(\mathbf{I}, \mathbf{I}, \mathbf{c}')^\dagger \stackrel{\text{a.s.}}{=} \mathbf{U}\mathbf{D}_c\mathbf{U}^\top (\mathbf{U}^\top)^\dagger \mathbf{D}_{c'}^{-1}\mathbf{U}^\dagger \stackrel{\text{a.s.}}{=} \mathbf{U}(\mathbf{D}_c\mathbf{D}_{c'}^{-1})\mathbf{U}^\dagger,$$

where $\mathbf{U} = [\boldsymbol{\mu}_1 | \dots | \boldsymbol{\mu}_K]$ are the linearly independent tensor components and $\mathbf{D}_c = \text{Diag}(\langle \boldsymbol{\mu}_1, \mathbf{c} \rangle, \dots, \langle \boldsymbol{\mu}_K, \mathbf{c} \rangle)$ is diagonal.

By linear independence of $\{\boldsymbol{\mu}_i\}_{i=1}^K$ and random choice of \mathbf{c} and \mathbf{c}' :

- 1 \mathbf{U} has rank K ;
- 2 \mathbf{D}_c and $\mathbf{D}_{c'}$ are invertible (a.s.);
- 3 Diagonal entries of $\mathbf{D}_c\mathbf{D}_{c'}^{-1}$ are distinct (a.s.);

So $\{\boldsymbol{\mu}_i\}_{i=1}^K$ are the eigenvectors of $\mathcal{T}(\mathbf{I}, \mathbf{I}, \mathbf{c})\mathcal{T}(\mathbf{I}, \mathbf{I}, \mathbf{c}')^\dagger$ with distinct non-zero eigenvalues.

Jennrich's algorithm is consistent

Error-tolerant algorithms for tensor decompositions

Moment Estimator: Empirical Moments

Moment Estimator: Empirical Moments

- Moments $\mathbb{E}_{\theta}[f(\mathbf{X})]$ are functions of model parameters θ
- Empirical Moments $\widehat{\mathbb{E}}[f(\mathbf{X})]$ are computed using iid samples $\{\mathbf{x}_i\}_{i=1}^n$ only

Moment Estimator: Empirical Moments

- Moments $\mathbb{E}_{\theta}[f(\mathbf{X})]$ are functions of model parameters θ
- Empirical Moments $\widehat{\mathbb{E}}[f(\mathbf{X})]$ are computed using iid samples $\{\mathbf{x}_i\}_{i=1}^n$ only

Example

- Third Order Moment: distribution of word triples

$$\mathbb{E}[\mathbf{x} \otimes \mathbf{x} \otimes \mathbf{x}] = \sum_h w_h \mathbf{a}_h \otimes \mathbf{a}_h \otimes \mathbf{a}_h$$

- Empirical Third Order Moment: co-occurrence frequency of word triples

$$\widehat{\mathbb{E}}[\mathbf{x} \otimes \mathbf{x} \otimes \mathbf{x}] = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \otimes \mathbf{x}_i \otimes \mathbf{x}_i$$

Moment Estimator: Empirical Moments

- Moments $\mathbb{E}_{\theta}[f(\mathbf{X})]$ are functions of model parameters θ
- Empirical Moments $\widehat{\mathbb{E}}[f(\mathbf{X})]$ are computed using iid samples $\{\mathbf{x}_i\}_{i=1}^n$ only

Example

- Third Order Moment: distribution of word triples

$$\mathbb{E}[\mathbf{x} \otimes \mathbf{x} \otimes \mathbf{x}] = \sum_h w_h \mathbf{a}_h \otimes \mathbf{a}_h \otimes \mathbf{a}_h$$

- Empirical Third Order Moment: co-occurrence frequency of word triples

$$\widehat{\mathbb{E}}[\mathbf{x} \otimes \mathbf{x} \otimes \mathbf{x}] = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \otimes \mathbf{x}_i \otimes \mathbf{x}_i$$

- Inevitably expect error of order $n^{-\frac{1}{2}}$ in some norm, e.g.,

$$\text{Operator norm: } \|\mathbb{E}[\mathbf{x} \otimes \mathbf{x} \otimes \mathbf{x}] - \widehat{\mathbb{E}}[\mathbf{x} \otimes \mathbf{x} \otimes \mathbf{x}]\| \lesssim n^{-\frac{1}{2}}$$

$$\text{where } \|\mathcal{T}\| := \sup_{\mathbf{x}, \mathbf{y}, \mathbf{z} \in S^{d-1}} \mathcal{T}(\mathbf{x}, \mathbf{y}, \mathbf{z})$$

$$\text{Frobenius norm: } \|\mathbb{E}[\mathbf{x} \otimes \mathbf{x} \otimes \mathbf{x}] - \widehat{\mathbb{E}}[\mathbf{x} \otimes \mathbf{x} \otimes \mathbf{x}]\|_F \lesssim n^{-\frac{1}{2}}$$

$$\text{where } \|\mathcal{T}\|_F := \sqrt{\sum_{i,j,k} T_{i,j,k}^2}$$

Stability of Jennrich's Algorithm

Recall Jennrich's algorithm

Given tensor $\mathcal{T} = \sum_{h=1}^K \boldsymbol{\mu}_h \otimes^3$ with linearly independent components $\{\boldsymbol{\mu}_h\}_{h=1}^K$, find the components (up to scaling).

Algorithm Jennrich's Algorithm

Require: Tensor $\mathcal{T} \in \mathbb{R}^{d \times d \times d}$

Ensure: Components $\{\hat{\boldsymbol{\mu}}_h\}_{h=1}^K \stackrel{\text{a.s.}}{=} \{\boldsymbol{\mu}_h\}_{h=1}^K$

- 1: Sample \mathbf{c} and \mathbf{c}' independently & uniformly at random from S^{d-1}
 - 2: Return $\{\hat{\boldsymbol{\mu}}_h\}_{h=1}^K \leftarrow$ eigenvectors of $(\mathcal{T}(\mathbf{I}, \mathbf{I}, \mathbf{c})\mathcal{T}(\mathbf{I}, \mathbf{I}, \mathbf{c}')^\dagger)$
-

Stability of Jennrich's Algorithm

Recall Jennrich's algorithm

Given tensor $\mathcal{T} = \sum_{h=1}^K \boldsymbol{\mu}_h \otimes^3$ with linearly independent components $\{\boldsymbol{\mu}_h\}_{h=1}^K$, find the components (up to scaling).

Algorithm Jennrich's Algorithm

Require: Tensor $\mathcal{T} \in \mathbb{R}^{d \times d \times d}$

Ensure: Components $\{\hat{\boldsymbol{\mu}}_h\}_{h=1}^K \stackrel{\text{a.s.}}{=} \{\boldsymbol{\mu}_h\}_{h=1}^K$

- 1: Sample \mathbf{c} and \mathbf{c}' independently & uniformly at random from S^{d-1}
 - 2: Return $\{\hat{\boldsymbol{\mu}}_h\}_{h=1}^K \leftarrow$ eigenvectors of $(\mathcal{T}(\mathbf{I}, \mathbf{I}, \mathbf{c})\mathcal{T}(\mathbf{I}, \mathbf{I}, \mathbf{c}')^\dagger)$
-

Challenge: Only have access to $\hat{\mathcal{T}}$ such that $\|\hat{\mathcal{T}} - \mathcal{T}\| \lesssim n^{-\frac{1}{2}}$

Stability of Jennrich's Algorithm

Recall Jennrich's algorithm

Given tensor $\mathcal{T} = \sum_{h=1}^K \boldsymbol{\mu}_h \otimes^3$ with linearly independent components $\{\boldsymbol{\mu}_h\}_{h=1}^K$, find the components (up to scaling).

Algorithm Jennrich's Algorithm

Require: Tensor $\hat{\mathcal{T}} \in \mathbb{R}^{d \times d \times d}$

Ensure: Components $\{\hat{\boldsymbol{\mu}}_h\}_{h=1}^K \stackrel{\text{a.s.}}{=} \{\boldsymbol{\mu}_h\}_{h=1}^K$?

- 1: Sample \mathbf{c} and \mathbf{c}' independently & uniformly at random from S^{d-1}
 - 2: Return $\{\hat{\boldsymbol{\mu}}_h\}_{h=1}^K \leftarrow$ eigenvectors of $(\hat{\mathcal{T}}(\mathbf{I}, \mathbf{I}, \mathbf{c})\hat{\mathcal{T}}(\mathbf{I}, \mathbf{I}, \mathbf{c}')^\dagger)$
-

Challenge: Only have access to $\hat{\mathcal{T}}$ such that $\|\hat{\mathcal{T}} - \mathcal{T}\| \lesssim n^{-\frac{1}{2}}$

Stability of Jennrich's Algorithm

Recall Jennrich's algorithm

Given tensor $\mathcal{T} = \sum_{h=1}^K \boldsymbol{\mu}_h \otimes^3$ with linearly independent components $\{\boldsymbol{\mu}_h\}_{h=1}^K$, find the components (up to scaling).

Algorithm Jennrich's Algorithm

Require: Tensor $\hat{\mathcal{T}} \in \mathbb{R}^{d \times d \times d}$

Ensure: Components $\{\hat{\boldsymbol{\mu}}_h\}_{h=1}^K \stackrel{\text{a.s.}}{=} \{\boldsymbol{\mu}_h\}_{h=1}^K$?

- 1: Sample \mathbf{c} and \mathbf{c}' independently & uniformly at random from S^{d-1}
 - 2: Return $\{\hat{\boldsymbol{\mu}}_h\}_{h=1}^K \leftarrow$ eigenvectors of $(\hat{\mathcal{T}}(\mathbf{I}, \mathbf{I}, \mathbf{c})\hat{\mathcal{T}}(\mathbf{I}, \mathbf{I}, \mathbf{c}')^\dagger)$
-

Stability of eigenvectors requires eigenvalue gaps

Stability of Jennrich's Algorithm

Recall Jennrich's algorithm

Given tensor $\mathcal{T} = \sum_{h=1}^K \boldsymbol{\mu}_h \otimes^3$ with linearly independent components $\{\boldsymbol{\mu}_h\}_{h=1}^K$, find the components (up to scaling).

Algorithm Jennrich's Algorithm

Require: Tensor $\hat{\mathcal{T}} \in \mathbb{R}^{d \times d \times d}$

Ensure: Components $\{\hat{\boldsymbol{\mu}}_h\}_{h=1}^K \stackrel{\text{a.s.}}{=} \{\boldsymbol{\mu}_h\}_{h=1}^K$?

- 1: Sample \mathbf{c} and \mathbf{c}' independently & uniformly at random from S^{d-1}
 - 2: Return $\{\hat{\boldsymbol{\mu}}_h\}_{h=1}^K \leftarrow$ eigenvectors of $(\hat{\mathcal{T}}(\mathbf{I}, \mathbf{I}, \mathbf{c})\hat{\mathcal{T}}(\mathbf{I}, \mathbf{I}, \mathbf{c}')^\dagger)$
-

Stability of eigenvectors requires eigenvalue gaps

- To ensure eigenvalue gaps for $\hat{\mathcal{T}}(\cdot, \cdot, \mathbf{c})\hat{\mathcal{T}}(\cdot, \cdot, \mathbf{c}')^\dagger$, $\|\hat{\mathcal{T}}(\cdot, \cdot, \mathbf{c})\hat{\mathcal{T}}(\cdot, \cdot, \mathbf{c}')^\dagger - \mathcal{T}(\cdot, \cdot, \mathbf{c})\mathcal{T}(\cdot, \cdot, \mathbf{c}')^\dagger\| \ll \Delta$ is needed.

Stability of Jennrich's Algorithm

Recall Jennrich's algorithm

Given tensor $\mathcal{T} = \sum_{h=1}^K \boldsymbol{\mu}_h \otimes^3$ with linearly independent components $\{\boldsymbol{\mu}_h\}_{h=1}^K$, find the components (up to scaling).

Algorithm Jennrich's Algorithm

Require: Tensor $\hat{\mathcal{T}} \in \mathbb{R}^{d \times d \times d}$

Ensure: Components $\{\hat{\boldsymbol{\mu}}_h\}_{h=1}^K \stackrel{\text{a.s.}}{=} \{\boldsymbol{\mu}_h\}_{h=1}^K$?

- 1: Sample \mathbf{c} and \mathbf{c}' independently & uniformly at random from S^{d-1}
 - 2: Return $\{\hat{\boldsymbol{\mu}}_h\}_{h=1}^K \leftarrow$ eigenvectors of $(\hat{\mathcal{T}}(\mathbf{I}, \mathbf{I}, \mathbf{c})\hat{\mathcal{T}}(\mathbf{I}, \mathbf{I}, \mathbf{c}')^\dagger)$
-

Stability of eigenvectors requires eigenvalue gaps

- To ensure eigenvalue gaps for $\hat{\mathcal{T}}(\cdot, \cdot, \mathbf{c})\hat{\mathcal{T}}(\cdot, \cdot, \mathbf{c}')^\dagger$, $\|\hat{\mathcal{T}}(\cdot, \cdot, \mathbf{c})\hat{\mathcal{T}}(\cdot, \cdot, \mathbf{c}')^\dagger - \mathcal{T}(\cdot, \cdot, \mathbf{c})\mathcal{T}(\cdot, \cdot, \mathbf{c}')^\dagger\| \ll \Delta$ is needed.
- Ultimately, $\|\hat{\mathcal{T}} - \mathcal{T}\|_F \ll \frac{1}{\text{poly } d}$ is required.

Stability of Jennrich's Algorithm

Recall Jennrich's algorithm

Given tensor $\mathcal{T} = \sum_{h=1}^K \boldsymbol{\mu}_h \otimes^3$ with linearly independent components $\{\boldsymbol{\mu}_h\}_{h=1}^K$, find the components (up to scaling).

Algorithm Jennrich's Algorithm

Require: Tensor $\hat{\mathcal{T}} \in \mathbb{R}^{d \times d \times d}$

Ensure: Components $\{\hat{\boldsymbol{\mu}}_h\}_{h=1}^K \stackrel{\text{a.s.}}{=} \{\boldsymbol{\mu}_h\}_{h=1}^K$?

- 1: Sample \mathbf{c} and \mathbf{c}' independently & uniformly at random from S^{d-1}
 - 2: Return $\{\hat{\boldsymbol{\mu}}_h\}_{h=1}^K \leftarrow$ eigenvectors of $(\hat{\mathcal{T}}(\mathbf{I}, \mathbf{I}, \mathbf{c})\hat{\mathcal{T}}(\mathbf{I}, \mathbf{I}, \mathbf{c}')^\dagger)$
-

Stability of eigenvectors requires eigenvalue gaps

- To ensure eigenvalue gaps for $\hat{\mathcal{T}}(\cdot, \cdot, \mathbf{c})\hat{\mathcal{T}}(\cdot, \cdot, \mathbf{c}')^\dagger$, $\|\hat{\mathcal{T}}(\cdot, \cdot, \mathbf{c})\hat{\mathcal{T}}(\cdot, \cdot, \mathbf{c}')^\dagger - \mathcal{T}(\cdot, \cdot, \mathbf{c})\mathcal{T}(\cdot, \cdot, \mathbf{c}')^\dagger\| \ll \Delta$ is needed.
- Ultimately, $\|\hat{\mathcal{T}} - \mathcal{T}\|_F \ll \frac{1}{\text{poly } d}$ is required. A different approach?

Initial Ideas

In many applications, we estimate moments of the form

$$\mathcal{M}_3 = \sum_{h=1}^K w_h \mathbf{a}_h \otimes^3,$$

where $\{\mathbf{a}_h\}_{h=1}^K$ are assumed to be linearly independent.

What if $\{\mathbf{a}_h\}_{h=1}^K$ has orthonormal columns?

Initial Ideas

In many applications, we estimate moments of the form

$$\mathcal{M}_3 = \sum_{h=1}^K w_h \mathbf{a}_h \otimes^3,$$

where $\{\mathbf{a}_h\}_{h=1}^K$ are assumed to be linearly independent.

What if $\{\mathbf{a}_h\}_{h=1}^K$ has orthonormal columns?

$$\mathcal{M}_3(I, \mathbf{a}_i, \mathbf{a}_i) = \sum_h w_h \langle \mathbf{a}_h, \mathbf{a}_i \rangle^2 \mathbf{a}_h = w_i \mathbf{a}_i, \forall i.$$

Initial Ideas

In many applications, we estimate moments of the form

$$\mathcal{M}_3 = \sum_{h=1}^K w_h \mathbf{a}_h \otimes^3,$$

where $\{\mathbf{a}_h\}_{h=1}^K$ are assumed to be linearly independent.

What if $\{\mathbf{a}_h\}_{h=1}^K$ has orthonormal columns?

$$\mathcal{M}_3(I, \mathbf{a}_i, \mathbf{a}_i) = \sum_h w_h \langle \mathbf{a}_h, \mathbf{a}_i \rangle^2 \mathbf{a}_h = w_i \mathbf{a}_i, \forall i.$$

- Analogous to matrix eigenvectors: $\mathbf{M}\mathbf{v} = \mathbf{M}(I, \mathbf{v}) = \lambda\mathbf{v}$.

Initial Ideas

In many applications, we estimate moments of the form

$$\mathcal{M}_3 = \sum_{h=1}^K w_h \mathbf{a}_h \otimes^3,$$

where $\{\mathbf{a}_h\}_{h=1}^K$ are assumed to be linearly independent.

What if $\{\mathbf{a}_h\}_{h=1}^K$ has orthonormal columns?

$$\mathcal{M}_3(I, \mathbf{a}_i, \mathbf{a}_i) = \sum_h w_h \langle \mathbf{a}_h, \mathbf{a}_i \rangle^2 \mathbf{a}_h = w_i \mathbf{a}_i, \forall i.$$

- Analogous to matrix eigenvectors: $\mathbf{M}\mathbf{v} = \mathbf{M}(I, \mathbf{v}) = \lambda\mathbf{v}$.
- **Define** orthonormal $\{\mathbf{a}_h\}_{h=1}^K$ as **eigenvectors** of tensor \mathcal{M}_3 .

Initial Ideas

In many applications, we estimate moments of the form

$$\mathcal{M}_3 = \sum_{h=1}^K w_h \mathbf{a}_h \otimes^3,$$

where $\{\mathbf{a}_h\}_{h=1}^K$ are assumed to be linearly independent.

What if $\{\mathbf{a}_h\}_{h=1}^K$ has orthonormal columns?

$$\mathcal{M}_3(I, \mathbf{a}_i, \mathbf{a}_i) = \sum_h w_h \langle \mathbf{a}_h, \mathbf{a}_i \rangle^2 \mathbf{a}_h = w_i \mathbf{a}_i, \forall i.$$

- Analogous to matrix eigenvectors: $M\mathbf{v} = M(I, \mathbf{v}) = \lambda\mathbf{v}$.
- **Define** orthonormal $\{\mathbf{a}_h\}_{h=1}^K$ as **eigenvectors** of tensor \mathcal{M}_3 .

Two Problems

- $\{\mathbf{a}_h\}_{h=1}^K$ is not orthogonal in general.
- How to find eigenvectors of a tensor?

Initial Ideas

In many applications, we estimate moments of the form

$$\mathcal{M}_3 = \sum_{h=1}^K w_h \mathbf{a}_h \otimes^3,$$

where $\{\mathbf{a}_h\}_{h=1}^K$ are assumed to be linearly independent.

What if $\{\mathbf{a}_h\}_{h=1}^K$ has orthonormal columns?

$$\mathcal{M}_3(I, \mathbf{a}_i, \mathbf{a}_i) = \sum_h w_h \langle \mathbf{a}_h, \mathbf{a}_i \rangle^2 \mathbf{a}_h = w_i \mathbf{a}_i, \forall i.$$

- Analogous to matrix eigenvectors: $M\mathbf{v} = M(I, \mathbf{v}) = \lambda\mathbf{v}$.
- **Define** orthonormal $\{\mathbf{a}_h\}_{h=1}^K$ as **eigenvectors** of tensor \mathcal{M}_3 .

Two Problems

- $\{\mathbf{a}_h\}_{h=1}^K$ **is not orthogonal in general.**
- How to find eigenvectors of a tensor?

Whitening is the process of finding a whitening matrix \mathbf{W} such that multi-linear operation (using \mathbf{W}) on \mathcal{M}_3 orthogonalize its components:

$$\begin{aligned}\mathcal{M}_3(\mathbf{W}, \mathbf{W}, \mathbf{W}) &= \sum_h w_h (\mathbf{W}^\top \mathbf{a}_h) \otimes^3 \\ &= \sum_h w_h \mathbf{v}_h \otimes^3, \quad \mathbf{v}_h \perp \mathbf{v}_{h'}, \quad \forall h \neq h'\end{aligned}$$

Whitening

Given

$$\mathcal{M}_3 = \sum_h w_h \mathbf{a}_h \otimes^3, \quad \mathcal{M}_2 = \sum_h w_h \mathbf{a}_h \otimes \mathbf{a}_h,$$

Whitening

Given

$$\mathcal{M}_3 = \sum_h w_h \mathbf{a}_h \otimes^3, \quad \mathcal{M}_2 = \sum_h w_h \mathbf{a}_h \otimes \mathbf{a}_h,$$

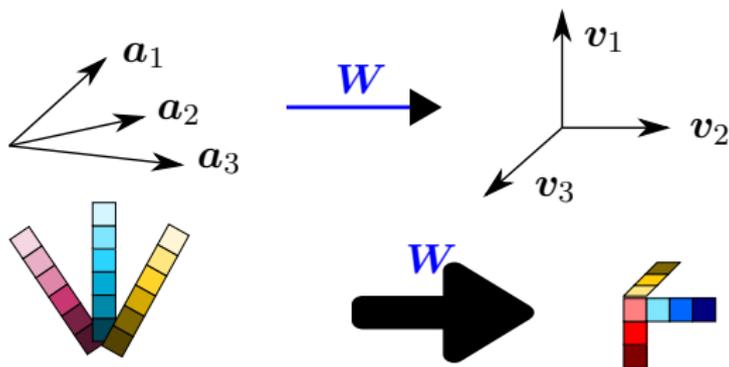
- Find whitening matrix \mathbf{W} s.t. $\mathbf{W}^\top \mathbf{a}_h = \mathbf{v}_h$ are orthogonal.

Whitening

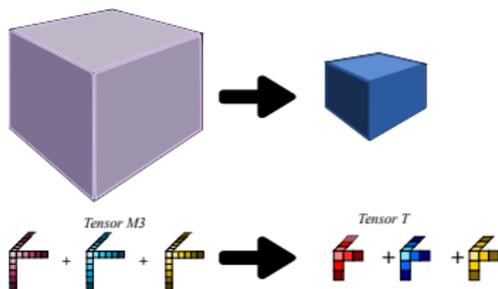
Given

$$\mathcal{M}_3 = \sum_h w_h \mathbf{a}_h \otimes \mathbf{a}_h, \quad \mathcal{M}_2 = \sum_h w_h \mathbf{a}_h \otimes \mathbf{a}_h,$$

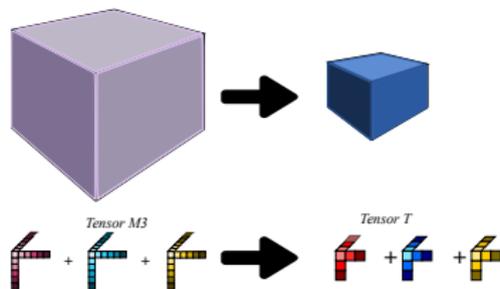
- Find whitening matrix \mathbf{W} s.t. $\mathbf{W}^\top \mathbf{a}_h = \mathbf{v}_h$ are orthogonal.
- When $\{\mathbf{a}_h\}_{h=1}^K \in \mathbb{R}^{d \times K}$ has full column rank, it is an invertible transformation.



Using Whitening to Obtain Orthogonal Tensor



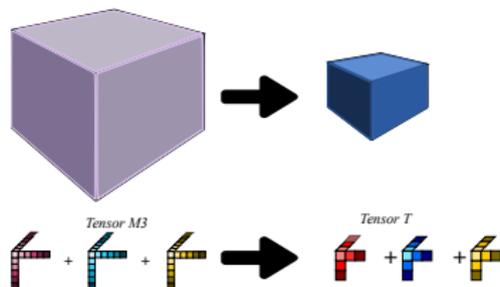
Using Whitening to Obtain Orthogonal Tensor



Multi-linear transform

- $\mathcal{T} = \mathcal{M}_3(\mathbf{W}, \mathbf{W}, \mathbf{W}) = \sum_h w_h (\mathbf{W}^\top \mathbf{a}_h)^{\otimes 3}$.

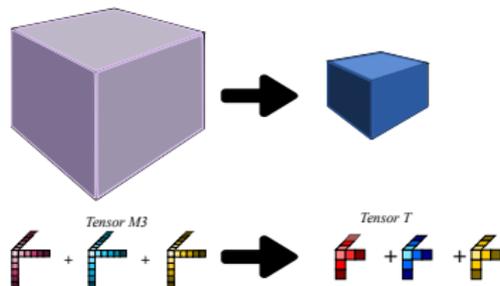
Using Whitening to Obtain Orthogonal Tensor



Multi-linear transform

- $\mathcal{T} = \mathcal{M}_3(\mathbf{W}, \mathbf{W}, \mathbf{W}) = \sum_h w_h (\mathbf{W}^\top \mathbf{a}_h)^{\otimes 3}$.
- $\mathcal{T} = \sum_{h \in [K]} w_h \cdot \mathbf{v}_h^{\otimes 3}$ has orthogonal components.

Using Whitening to Obtain Orthogonal Tensor



Multi-linear transform

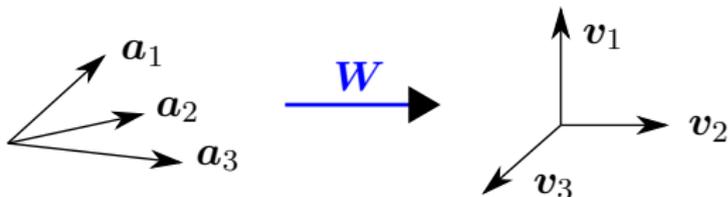
- $\mathcal{T} = \mathcal{M}_3(\mathbf{W}, \mathbf{W}, \mathbf{W}) = \sum_h w_h (\mathbf{W}^\top \mathbf{a}_h)^{\otimes 3}$.
- $\mathcal{T} = \sum_{h \in [K]} w_h \cdot \mathbf{v}_h^{\otimes 3}$ has orthogonal components.
- Dimensionality reduction when $K \ll d$, as $\mathcal{M}_3 \in \mathbb{R}^{d \times d \times d}$ and $\mathcal{T} \in \mathbb{R}^{K \times K \times K}$.

How to Find Whitening Matrix?

Given

$$\mathcal{M}_3 = \sum_h w_h \mathbf{a}_h \otimes^3, \quad \mathcal{M}_2 = \sum_h w_h \mathbf{a}_h \otimes \mathbf{a}_h,$$

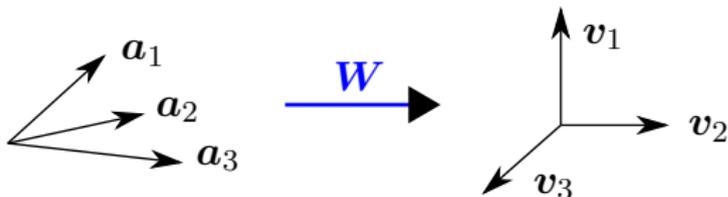
Goal: W such that



How to Find Whitening Matrix?

Given

$$\mathcal{M}_3 = \sum_h w_h \mathbf{a}_h \otimes^3, \quad \mathcal{M}_2 = \sum_h w_h \mathbf{a}_h \otimes \mathbf{a}_h,$$



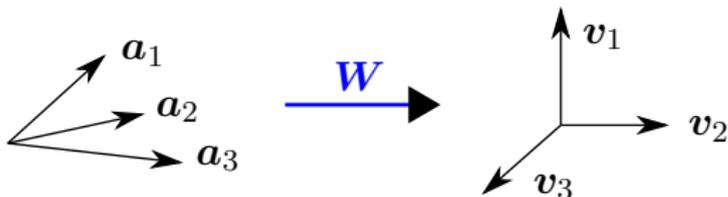
Goal: W such that

- Use pairwise moments \mathcal{M}_2 to find W s.t. $W^\top \mathcal{M}_2 W = I$.

How to Find Whitening Matrix?

Given

$$\mathcal{M}_3 = \sum_h w_h \mathbf{a}_h \otimes^3, \quad M_2 = \sum_h w_h \mathbf{a}_h \otimes \mathbf{a}_h,$$



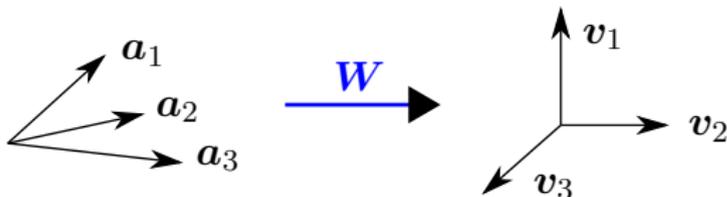
Goal: W such that

- Use pairwise moments M_2 to find W s.t. $W^\top M_2 W = I$.
- $W = U \text{Diag}(\tilde{\lambda}^{-1/2})$, where Eigen-decomposition $M_2 = U \text{Diag}(\tilde{\lambda}) U^\top$.

How to Find Whitening Matrix?

Given

$$\mathcal{M}_3 = \sum_h w_h \mathbf{a}_h \otimes^3, \quad M_2 = \sum_h w_h \mathbf{a}_h \otimes \mathbf{a}_h,$$



Goal: W such that

- Use pairwise moments M_2 to find W s.t. $W^\top M_2 W = I$.
- $W = U \text{Diag}(\tilde{\lambda}^{-1/2})$, where Eigen-decomposition $M_2 = U \text{Diag}(\tilde{\lambda}) U^\top$.
- $V := W^\top A \text{Diag}(w)^{1/2}$ is an orthogonal matrix.

$$\begin{aligned} \mathcal{T} = \mathcal{M}_3(W, W, W) &= \sum_h w_h^{-1/2} (W^\top \mathbf{a}_h \sqrt{w_h}) \otimes^3 \\ &= \sum_h \lambda_h \mathbf{v}_h \otimes^3, \quad \lambda_h := w_h^{-1/2}. \end{aligned}$$

\mathcal{T} is an orthogonal tensor.

Initial Ideas

In many applications, we estimate moments of the form

$$\mathcal{M}_3 = \sum_h w_h \mathbf{a}_h \otimes^3,$$

where $\{\mathbf{a}_h\}_{h=1}^K$ are assumed to be linearly independent.

What if $\{\mathbf{a}_h\}_{h=1}^K$ has orthonormal columns?

$$\mathcal{M}_3(I, \mathbf{a}_i, \mathbf{a}_i) = \sum_h w_h \langle \mathbf{a}_h, \mathbf{a}_i \rangle^2 \mathbf{a}_h = w_i \mathbf{a}_i, \forall i.$$

- Analogous to matrix eigenvectors: $M\mathbf{v} = M(I, \mathbf{v}) = \lambda\mathbf{v}$.
- Define orthonormal $\{\mathbf{a}_h\}_{h=1}^K$ as eigenvectors of tensor \mathcal{M}_3 .

Two Problems

- $\{\mathbf{a}_h\}_{h=1}^K$ is not orthogonal in general.
- **How to find eigenvectors of a tensor?**

Review: Orthogonal Matrix Eigen Decomposition

Task: Given matrix $\mathbf{M} = \sum_{h=1}^K \lambda_h \mathbf{v}_h \otimes \mathbf{v}_h$ with orthonormal components $\{\mathbf{v}_h\}_{h=1}^K$ ($\mathbf{v}_h \perp \mathbf{v}_{h'}, \forall h \neq h'$), find the components/eigenvectors.

Review: Orthogonal Matrix Eigen Decomposition

Task: Given matrix $M = \sum_{h=1}^K \lambda_h \mathbf{v}_h \otimes \mathbf{v}_h$ with orthonormal components $\{\mathbf{v}_h\}_{h=1}^K$ ($\mathbf{v}_h \perp \mathbf{v}_{h'}, \forall h \neq h'$), find the components/eigenvectors.

Properties of Matrix Eigenvectors

- Fixed point: linear transform $M(\mathbf{I}, \mathbf{v}_i) = \sum_h \lambda_h \langle \mathbf{v}_i, \mathbf{v}_h \rangle \mathbf{v}_h = \lambda_i \mathbf{v}_i$

Review: Orthogonal Matrix Eigen Decomposition

Task: Given matrix $M = \sum_{h=1}^K \lambda_h \mathbf{v}_h \otimes \mathbf{v}_h$ with orthonormal components $\{\mathbf{v}_h\}_{h=1}^K$ ($\mathbf{v}_h \perp \mathbf{v}_{h'}, \forall h \neq h'$), find the components/eigenvectors.

Properties of Matrix Eigenvectors

- Fixed point: linear transform $M(\mathbf{I}, \mathbf{v}_i) = \sum_h \lambda_h \langle \mathbf{v}_i, \mathbf{v}_h \rangle \mathbf{v}_h = \lambda_i \mathbf{v}_i$

Intuitions for Matrix Power Method

Review: Orthogonal Matrix Eigen Decomposition

Task: Given matrix $M = \sum_{h=1}^K \lambda_h \mathbf{v}_h \otimes \mathbf{v}_h$ with orthonormal components $\{\mathbf{v}_h\}_{h=1}^K$ ($\mathbf{v}_h \perp \mathbf{v}_{h'}, \forall h \neq h'$), find the components/eigenvectors.

Properties of Matrix Eigenvectors

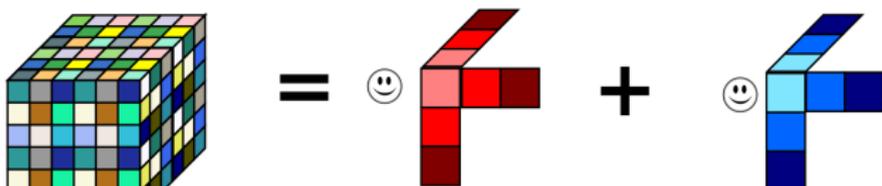
- Fixed point: linear transform $M(\mathbf{I}, \mathbf{v}_i) = \sum_h \lambda_h \langle \mathbf{v}_i, \mathbf{v}_h \rangle \mathbf{v}_h = \lambda_i \mathbf{v}_i$

Intuitions for Matrix Power Method

Linear transform on eigenvectors $\{\mathbf{v}_h\}_{h=1}^K$ preserve direction

Orthogonal Tensor Eigen Decomposition

Task: Given tensor $\mathcal{T} = \sum_{h=1}^K \lambda_h \mathbf{v}_h \otimes^3$ with orthonormal components $\{\mathbf{v}_h\}_{h=1}^K$ ($\mathbf{v}_h \perp \mathbf{v}_{h'}, \forall h \neq h'$), find the components/eigenvectors.



$Tensor = u_1 \otimes u_1 \otimes u_1 + u_2 \otimes u_2 \otimes u_2$

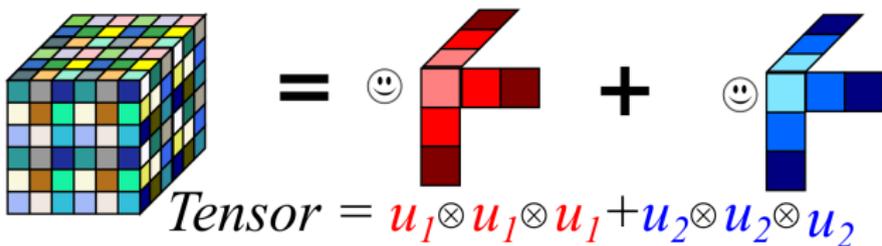
Orthogonal Tensor Eigen Decomposition

Task: Given tensor $\mathcal{T} = \sum_{h=1}^K \lambda_h \mathbf{v}_h \otimes^3$ with orthonormal components $\{\mathbf{v}_h\}_{h=1}^K$ ($\mathbf{v}_h \perp \mathbf{v}_{h'}, \forall h \neq h'$), find the components/eigenvectors.

Properties of Tensor Eigenvectors

- Fixed point: bi-linear transform

$$\mathcal{T}(\mathbf{I}, \mathbf{v}_i, \mathbf{v}_i) = \sum_h \lambda_h \langle \mathbf{v}_i, \mathbf{v}_h \rangle^2 \mathbf{v}_h = \lambda_i \mathbf{v}_i$$



$Tensor = u_1 \otimes u_1 \otimes u_1 + u_2 \otimes u_2 \otimes u_2$

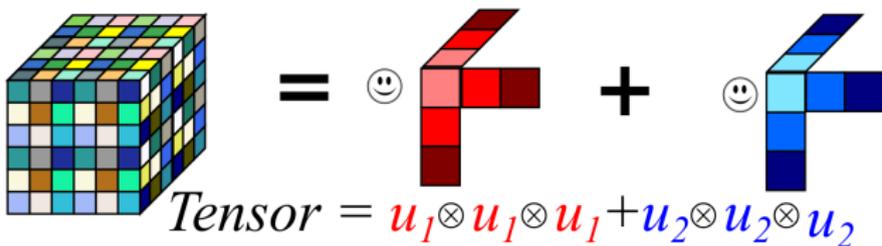
Orthogonal Tensor Eigen Decomposition

Task: Given tensor $\mathcal{T} = \sum_{h=1}^K \lambda_h \mathbf{v}_h \otimes^3$ with orthonormal components $\{\mathbf{v}_h\}_{h=1}^K$ ($\mathbf{v}_h \perp \mathbf{v}_{h'}, \forall h \neq h'$), find the components/eigenvectors.

Properties of Tensor Eigenvectors

- Fixed point: bi-linear transform

$$\mathcal{T}(\mathbf{I}, \mathbf{v}_i, \mathbf{v}_i) = \sum_h \lambda_h \langle \mathbf{v}_i, \mathbf{v}_h \rangle^2 \mathbf{v}_h = \lambda_i \mathbf{v}_i$$



$Tensor = u_1 \otimes u_1 \otimes u_1 + u_2 \otimes u_2 \otimes u_2$

Intuitions for Tensor Power Method

Orthogonal Tensor Eigen Decomposition

Task: Given tensor $\mathcal{T} = \sum_{h=1}^K \lambda_h \mathbf{v}_h \otimes^3$ with orthonormal components $\{\mathbf{v}_h\}_{h=1}^K$ ($\mathbf{v}_h \perp \mathbf{v}_{h'}, \forall h \neq h'$), find the components/eigenvectors.

Properties of Tensor Eigenvectors

- Fixed point: bi-linear transform

$$\mathcal{T}(\mathbf{I}, \mathbf{v}_i, \mathbf{v}_i) = \sum_h \lambda_h \langle \mathbf{v}_i, \mathbf{v}_h \rangle^2 \mathbf{v}_h = \lambda_i \mathbf{v}_i$$

$Tensor = \mathbf{u}_1 \otimes \mathbf{u}_1 \otimes \mathbf{u}_1 + \mathbf{u}_2 \otimes \mathbf{u}_2 \otimes \mathbf{u}_2$

Intuitions for Tensor Power Method

Bilinear transform on eigenvectors $\{\mathbf{v}_h\}_{h=1}^K$ preserve direction

Orthogonal Matrix Eigen Decomposition

Task: Given matrix $M = \sum_{h=1}^K \lambda_h \mathbf{v}_h \otimes^2$ with orthonormal components $\{\mathbf{v}_h\}_{h=1}^K$ ($\mathbf{v}_h \perp \mathbf{v}_{h'}, \forall h \neq h'$), find the components/eigenvectors.

Algorithm Matrix Power Method

Require: Matrix $M \in \mathbb{R}^{K \times K}$

Ensure: Components $\{\hat{\mathbf{v}}_h\}_{h=1}^K \stackrel{\text{w.h.p.}}{=} \{\mathbf{v}_h\}_{h=1}^K$

- 1: **for** $h = 1 : K$ **do**
 - 2: Sample \mathbf{u}_0 uniformly at random from S^{K-1}
 - 3: **for** $i = 1 : T$ **do**
 - 4: $\mathbf{u}_i \leftarrow \frac{M(\mathbf{I}, \mathbf{u}_{i-1})}{\|M(\mathbf{I}, \mathbf{u}_{i-1})\|}$
 - 5: **end for**
 - 6: $\hat{\mathbf{v}}_h \leftarrow \mathbf{u}_T, \hat{\lambda}_h \leftarrow M(\hat{\mathbf{v}}_h, \hat{\mathbf{v}}_h)$
 - 7: Deflate $M \leftarrow M - \hat{\lambda}_h \hat{\mathbf{v}}_h \otimes^2$
 - 8: **end for**
-

Orthogonal Matrix Eigen Decomposition

Task: Given matrix $M = \sum_{h=1}^K \lambda_h \mathbf{v}_h \otimes \mathbf{v}_h^T$ with orthonormal components $\{\mathbf{v}_h\}_{h=1}^K$ ($\mathbf{v}_h \perp \mathbf{v}_{h'}, \forall h \neq h'$), find the components/eigenvectors.

Algorithm Matrix Power Method

Require: Matrix $M \in \mathbb{R}^{K \times K}$

Ensure: Components $\{\hat{\mathbf{v}}_h\}_{h=1}^K \stackrel{\text{w.h.p.}}{=} \{\mathbf{v}_h\}_{h=1}^K$

- 1: **for** $h = 1 : K$ **do**
 - 2: Sample \mathbf{u}_0 uniformly at random from S^{K-1}
 - 3: **for** $i = 1 : T$ **do**
 - 4: $\mathbf{u}_i \leftarrow \frac{M(\mathbf{I}, \mathbf{u}_{i-1})}{\|M(\mathbf{I}, \mathbf{u}_{i-1})\|}$
 - 5: **end for**
 - 6: $\hat{\mathbf{v}}_h \leftarrow \mathbf{u}_T, \hat{\lambda}_h \leftarrow M(\hat{\mathbf{v}}_h, \hat{\mathbf{v}}_h)$
 - 7: Deflate $M \leftarrow M - \hat{\lambda}_h \hat{\mathbf{v}}_h \otimes \hat{\mathbf{v}}_h^T$
 - 8: **end for**
-

Consistency of Matrix Power Method?

Is there convergence? $\{\hat{\mathbf{v}}_h\}_{h=1}^K \equiv \{\mathbf{v}_h\}_{h=1}^K$ w.h.p.?

Orthogonal Matrix Eigen Decomposition

Task: Given matrix $M = \sum_{h=1}^K \lambda_h \mathbf{v}_h \otimes \mathbf{v}_h^T$ with orthonormal components $\{\mathbf{v}_h\}_{h=1}^K$ ($\mathbf{v}_h \perp \mathbf{v}_{h'}, \forall h \neq h'$), find the components/eigenvectors.

Algorithm Matrix Power Method

Require: Matrix $M \in \mathbb{R}^{K \times K}$

Ensure: Components $\{\hat{\mathbf{v}}_h\}_{h=1}^K \stackrel{\text{w.h.p.}}{=} \{\mathbf{v}_h\}_{h=1}^K$

- 1: **for** $h = 1 : K$ **do**
 - 2: Sample \mathbf{u}_0 uniformly at random from S^{K-1}
 - 3: **for** $i = 1 : T$ **do**
 - 4: $\mathbf{u}_i \leftarrow \frac{M(\mathbf{I}, \mathbf{u}_{i-1})}{\|M(\mathbf{I}, \mathbf{u}_{i-1})\|}$
 - 5: **end for**
 - 6: $\hat{\mathbf{v}}_h \leftarrow \mathbf{u}_T, \hat{\lambda}_h \leftarrow M(\hat{\mathbf{v}}_h, \hat{\mathbf{v}}_h)$
 - 7: Deflate $M \leftarrow M - \hat{\lambda}_h \hat{\mathbf{v}}_h \otimes \hat{\mathbf{v}}_h^T$
 - 8: **end for**
-

Consistency of Matrix Power Method?

Is there convergence? $\{\hat{\mathbf{v}}_h\}_{h=1}^K \equiv \{\mathbf{v}_h\}_{h=1}^K$ w.h.p.?

Does the convergence depend on initialization?

Orthogonal Tensor Eigen Decomposition

Task: Given tensor $\mathcal{T} = \sum_{h=1}^K \lambda_h \mathbf{v}_h \otimes^3$ with orthonormal components $\{\mathbf{v}_h\}_{h=1}^K$ ($\mathbf{v}_h \perp \mathbf{v}_{h'}, \forall h \neq h'$), find the components/eigenvectors.

Algorithm Tensor Power Method

Require: Tensor $\mathcal{T} \in \mathbb{R}^{K \times K \times K}$

Ensure: Components $\{\hat{\mathbf{v}}_h\}_{h=1}^K \stackrel{\text{w.h.p.}}{=} \{\mathbf{v}_h\}_{h=1}^K$

1: **for** $h = 1 : K$ **do**

2: Sample \mathbf{u}_0 uniformly at random from S^{K-1}

3: **for** $i = 1 : T$ **do**

4: $\mathbf{u}_i \leftarrow \frac{\mathcal{T}(\mathbf{I}, \mathbf{u}_{i-1}, \mathbf{u}_{i-1})}{\|\mathcal{T}(\mathbf{I}, \mathbf{u}_{i-1}, \mathbf{u}_{i-1})\|}$

5: **end for**

6: $\hat{\mathbf{v}}_h \leftarrow \mathbf{u}_T, \hat{\lambda}_h \leftarrow \mathcal{T}(\hat{\mathbf{v}}_h, \hat{\mathbf{v}}_h, \hat{\mathbf{v}}_h)$

7: Deflate $\mathcal{T} \leftarrow \mathcal{T} - \hat{\lambda}_h \hat{\mathbf{v}}_h \otimes^3$

8: **end for**

Orthogonal Tensor Eigen Decomposition

Task: Given tensor $\mathcal{T} = \sum_{h=1}^K \lambda_h \mathbf{v}_h \otimes^3$ with orthonormal components $\{\mathbf{v}_h\}_{h=1}^K$ ($\mathbf{v}_h \perp \mathbf{v}_{h'}, \forall h \neq h'$), find the components/eigenvectors.

Algorithm Tensor Power Method

Require: Tensor $\mathcal{T} \in \mathbb{R}^{K \times K \times K}$

Ensure: Components $\{\hat{\mathbf{v}}_h\}_{h=1}^K \stackrel{\text{w.h.p.}}{=} \{\mathbf{v}_h\}_{h=1}^K$

- 1: **for** $h = 1 : K$ **do**
 - 2: Sample \mathbf{u}_0 uniformly at random from S^{K-1}
 - 3: **for** $i = 1 : T$ **do**
 - 4: $\mathbf{u}_i \leftarrow \frac{\mathcal{T}(\mathbf{I}, \mathbf{u}_{i-1}, \mathbf{u}_{i-1})}{\|\mathcal{T}(\mathbf{I}, \mathbf{u}_{i-1}, \mathbf{u}_{i-1})\|}$
 - 5: **end for**
 - 6: $\hat{\mathbf{v}}_h \leftarrow \mathbf{u}_T, \hat{\lambda}_h \leftarrow \mathcal{T}(\hat{\mathbf{v}}_h, \hat{\mathbf{v}}_h, \hat{\mathbf{v}}_h)$
 - 7: Deflate $\mathcal{T} \leftarrow \mathcal{T} - \hat{\lambda}_h \hat{\mathbf{v}}_h \otimes^3$
 - 8: **end for**
-

Consistency of Tensor Power Method?

Is there convergence? $\{\hat{\mathbf{v}}_h\}_{h=1}^K \equiv \{\mathbf{v}_h\}_{h=1}^K$ w.h.p.?

Does the convergence depend on initialization?

Analysis of Consistency of Matrix Power Method

- Order eigenvectors $\{\mathbf{v}_h\}_{h=1}^K$ such that corresponding eigenvalues satisfy $\lambda_1 \geq \lambda_2 \dots \geq \lambda_K$.
- Project initial point \mathbf{u}_0 onto eigenvectors $\{\mathbf{v}_h\}_{h=1}^K$

$$c_h = \langle \mathbf{u}_0, \mathbf{v}_h \rangle, \quad \forall h$$

Convergence properties

- **Unique (identifiable)** i.f.f. $\{\lambda_h\}_{h=1}^K$ are distinct.
- If gap $\frac{\lambda_2}{\lambda_1} < 1$ and $c_1 \neq 0$, matrix power method converges to \mathbf{v}_1 .
- Converges **linearly** to \mathbf{v}_1 assuming gap $\lambda_2/\lambda_1 < 1$.
 - ▶ Linear transform permits $M(\mathbf{I}, \mathbf{u}_0) = \sum_h \lambda_h (\mathbf{v}_h^\top \mathbf{u}_0) \mathbf{v}_h = \sum_h \lambda_h c_h \mathbf{v}_h$, i.e., projection in \mathbf{v}_h direction is scaled by λ_h .
 - ▶ In t iterations, $\frac{(\mathbf{v}_1^\top \mathbf{v})^2}{\sum_i (\mathbf{v}_i^\top \mathbf{v})^2} \geq 1 - K \left(\frac{\lambda_2}{\lambda_1} \right)^{2t}$.

Analysis of Consistency of Tensor Power Method

- Project initial point \mathbf{u}_0 onto eigenvectors $c_h = \langle \mathbf{u}_0, \mathbf{v}_h \rangle$, $\forall h$.
- Order eigenvectors $\{\mathbf{v}_h\}_{h=1}^K$ such that

$$\lambda_1 |c_1| > \lambda_2 |c_2| \geq \dots \geq \lambda_K |c_K|.$$

Convergence properties

- **Identifiable** i.f.f. $\{\lambda_h |c_h|\}_{h=1}^K$ are distinct. **Initialization dependent.**
- If $\frac{\lambda_2 |c_2|}{\lambda_1 |c_1|} < 1$ and $\lambda_1 |c_1| \neq 0$, tensor power method converges to \mathbf{v}_1 .
Note \mathbf{v}_1 is NOT necessarily the largest eigenvector.
- Converges **quadratically** to \mathbf{v}_1 assuming gap $\frac{\lambda_2 |c_2|}{\lambda_1 |c_1|} < 1$.
 - ▶ Bi-linear transform permits $\mathcal{T}(I, \mathbf{u}_0, \mathbf{u}_0) = \sum_h \lambda_h (\mathbf{v}_h^\top \mathbf{u}_0)^2 \mathbf{v}_h = \sum_h \lambda_h c_h^2 \mathbf{v}_h$
i.e., projection in \mathbf{v}_h direction is **squared** then scaled by λ_h .
 - ▶ In t iterations, $\frac{(\mathbf{v}_1^\top \mathbf{v})^2}{\sum_i (\mathbf{v}_i^\top \mathbf{v})^2} \geq 1 - k \left(\frac{\lambda_1}{\max_{i \neq 1} \lambda_i} \right)^2 \left| \frac{v_2 c_2}{v_1 c_1} \right|^{2^{t+1}}$.

Matrix vs. tensor power iteration

Matrix power iteration:

Tensor power iteration:

Matrix vs. tensor power iteration

Matrix power iteration:

- 1 Requires gap between largest and second-largest eigenvalue.
Property of the matrix only.

Tensor power iteration:

- 1 Requires gap between largest and second-largest $\lambda_h |c_h|$.
Property of the tensor and initialization u_0 .

Matrix vs. tensor power iteration

Matrix power iteration:

- 1 Requires gap between largest and second-largest eigenvalue.
Property of the matrix only.
- 2 Converges to top eigenvector.

Tensor power iteration:

- 1 Requires gap between largest and second-largest $\lambda_h |c_h|$.
Property of the tensor and initialization u_0 .
- 2 Converges to v_i which is the largest $v_h |c_h|$. Not necessarily the largest eigenvector.

Matrix vs. tensor power iteration

Matrix power iteration:

- 1 Requires gap between largest and second-largest eigenvalue.
Property of the matrix only.
- 2 Converges to top eigenvector.
- 3 Linear convergence. Need $O(\log(1/\epsilon))$ iterations.

Tensor power iteration:

- 1 Requires gap between largest and second-largest $\lambda_h |c_h|$.
Property of the tensor and initialization u_0 .
- 2 Converges to v_i which is the largest $v_h |c_h|$. Not necessarily the largest eigenvector.
- 3 Quadratic convergence. Need $O(\log \log(1/\epsilon))$ iterations.

Spurious Eigenvectors for Tensor Eigen Decomposition

$$\mathcal{T} = \sum_{h \in [K]} \lambda_h \mathbf{v}_h \otimes^3.$$

Characterization of eigenvectors: $\mathcal{T}(\mathbf{I}, \mathbf{v}, \mathbf{v}) = \lambda \mathbf{v}$?

- $\{\mathbf{v}_h\}_{h=1}^K$ are eigenvectors as $\mathcal{T}(\mathbf{I}, \mathbf{v}_h, \mathbf{v}_h) = \lambda_h \mathbf{v}_h$.

Spurious Eigenvectors for Tensor Eigen Decomposition

$$\mathcal{T} = \sum_{h \in [K]} \lambda_h \mathbf{v}_h \otimes^3.$$

Characterization of eigenvectors: $\mathcal{T}(\mathbf{I}, \mathbf{v}, \mathbf{v}) = \lambda \mathbf{v}$

- $\{\mathbf{v}_h\}_{h=1}^K$ are eigenvectors as $\mathcal{T}(\mathbf{I}, \mathbf{v}_h, \mathbf{v}_h) = \lambda_h \mathbf{v}_h$.
- **Bad news:** There can be other eigenvectors (unlike matrix case).

E.g., when $\{\lambda_h\}_{h=1}^K \equiv 1$

$$\mathbf{v} = \frac{\mathbf{v}_1 + \mathbf{v}_2}{\sqrt{2}} \text{ satisfies } \mathcal{T}(\mathbf{I}, \mathbf{v}, \mathbf{v}) = \frac{1}{\sqrt{2}} \mathbf{v}.$$

Spurious Eigenvectors for Tensor Eigen Decomposition

$$\mathcal{T} = \sum_{h \in [K]} \lambda_h \mathbf{v}_h \otimes^3.$$

Characterization of eigenvectors: $\mathcal{T}(\mathbf{I}, \mathbf{v}, \mathbf{v}) = \lambda \mathbf{v}$

- $\{\mathbf{v}_h\}_{h=1}^K$ are eigenvectors as $\mathcal{T}(\mathbf{I}, \mathbf{v}_h, \mathbf{v}_h) = \lambda_h \mathbf{v}_h$.
- **Bad news:** There can be other eigenvectors (unlike matrix case).

E.g., when $\{\lambda_h\}_{h=1}^K \equiv 1$

$$\mathbf{v} = \frac{\mathbf{v}_1 + \mathbf{v}_2}{\sqrt{2}} \text{ satisfies } \mathcal{T}(\mathbf{I}, \mathbf{v}, \mathbf{v}) = \frac{1}{\sqrt{2}} \mathbf{v}.$$

How do we avoid spurious solutions (not components $\{\mathbf{v}_h\}_{h=1}^K$)?

Spurious Eigenvectors for Tensor Eigen Decomposition

$$\mathcal{T} = \sum_{h \in [K]} \lambda_h \mathbf{v}_h \otimes^3.$$

Characterization of eigenvectors: $\mathcal{T}(\mathbf{I}, \mathbf{v}, \mathbf{v}) = \lambda \mathbf{v}$

- $\{\mathbf{v}_h\}_{h=1}^K$ are eigenvectors as $\mathcal{T}(\mathbf{I}, \mathbf{v}_h, \mathbf{v}_h) = \lambda_h \mathbf{v}_h$.
- **Bad news:** There can be other eigenvectors (unlike matrix case).

E.g., when $\{\lambda_h\}_{h=1}^K \equiv 1$

$$\mathbf{v} = \frac{\mathbf{v}_1 + \mathbf{v}_2}{\sqrt{2}} \text{ satisfies } \mathcal{T}(\mathbf{I}, \mathbf{v}, \mathbf{v}) = \frac{1}{\sqrt{2}} \mathbf{v}.$$

How do we avoid spurious solutions (not components $\{\mathbf{v}_h\}_{h=1}^K$)?

Optimization viewpoint of tensor Eigen decomposition will help.

Spurious Eigenvectors for Tensor Eigen Decomposition

$$\mathcal{T} = \sum_{h \in [K]} \lambda_h \mathbf{v}_h \otimes^3.$$

Characterization of eigenvectors: $\mathcal{T}(\mathbf{I}, \mathbf{v}, \mathbf{v}) = \lambda \mathbf{v}$?

- $\{\mathbf{v}_h\}_{h=1}^K$ are eigenvectors as $\mathcal{T}(\mathbf{I}, \mathbf{v}_h, \mathbf{v}_h) = \lambda_h \mathbf{v}_h$.
- **Bad news:** There can be other eigenvectors (unlike matrix case).

E.g., when $\{\lambda_h\}_{h=1}^K \equiv 1$

$$\mathbf{v} = \frac{\mathbf{v}_1 + \mathbf{v}_2}{\sqrt{2}} \text{ satisfies } \mathcal{T}(\mathbf{I}, \mathbf{v}, \mathbf{v}) = \frac{1}{\sqrt{2}} \mathbf{v}.$$

How do we avoid spurious solutions (not components $\{\mathbf{v}_h\}_{h=1}^K$)?

Optimization viewpoint of tensor Eigen decomposition will help.

All spurious eigenvectors are saddle points.

Optimization Viewpoint of Matrix/Tensor Eigen Decomposition

Optimization Viewpoint of Matrix/Tensor Eigen Decomposition

Optimization Problem

Matrix: $\max_v M(v, v)$ s.t. $\|v\| = 1$.

- Lagrangian:

$$L(v, \lambda) := M(v, v) - \lambda(v^\top v - 1).$$

Tensor: $\max_v T(v, v, v)$ s.t. $\|v\| = 1$.

- Lagrangian:

$$L(v, \lambda) := T(v, v, v) - 1.5\lambda(v^\top v - 1).$$

Optimization Viewpoint of Matrix/Tensor Eigen Decomposition

Optimization Problem

Matrix: $\max_v M(v, v)$ s.t. $\|v\| = 1$.

- Lagrangian:

$$L(v, \lambda) := M(v, v) - \lambda(v^\top v - 1).$$

Tensor: $\max_v T(v, v, v)$ s.t. $\|v\| = 1$.

- Lagrangian:

$$L(v, \lambda) := T(v, v, v) - 1.5\lambda(v^\top v - 1).$$

Non-convex: stationary points = {global optima, local optima, saddle point}

Optimization Viewpoint of Matrix/Tensor Eigen Decomposition

Optimization Problem

Matrix: $\max_v M(v, v)$ s.t. $\|v\| = 1$.

- Lagrangian:

$$L(v, \lambda) := M(v, v) - \lambda(v^\top v - 1).$$

Tensor: $\max_v T(v, v, v)$ s.t. $\|v\| = 1$.

- Lagrangian:

$$L(v, \lambda) := T(v, v, v) - 1.5\lambda(v^\top v - 1).$$

Non-convex: stationary points = {global optima, local optima, saddle point}

Stationary Points: first derivative $\nabla L(v, \lambda) = 0$

$$\nabla L(v, \lambda) = 2(M(I, v) - \lambda v) = 0$$

- Eigenvectors are stationary points.
- Power method $v \leftarrow \frac{M(I, v)}{\|M(I, v)\|}$ is a version of **gradient ascent**.

$$\nabla L(v, \lambda) = 3(T(I, v, v) - \lambda v) = 0$$

- Eigenvectors are stationary points.
- Power method $v \leftarrow \frac{T(I, v, v)}{\|T(I, v, v)\|}$ is a version of **gradient ascent**.

Optimization Viewpoint of Matrix/Tensor Eigen Decomposition

Optimization Problem

Matrix: $\max_v M(v, v)$ s.t. $\|v\| = 1$.

- Lagrangian:

$$L(v, \lambda) := M(v, v) - \lambda(v^\top v - 1).$$

Tensor: $\max_v T(v, v, v)$ s.t. $\|v\| = 1$.

- Lagrangian:

$$L(v, \lambda) := T(v, v, v) - 1.5\lambda(v^\top v - 1).$$

Non-convex: stationary points = {global optima, local optima, saddle point}

Stationary Points: first derivative $\nabla L(v, \lambda) = 0$

$$\nabla L(v, \lambda) = 2(M(I, v) - \lambda v) = 0$$

- Eigenvectors are stationary points.
- Power method $v \leftarrow \frac{M(I, v)}{\|M(I, v)\|}$ is a version of **gradient ascent**.

$$\nabla L(v, \lambda) = 3(T(I, v, v) - \lambda v) = 0$$

- Eigenvectors are stationary points.
- Power method $v \leftarrow \frac{T(I, v, v)}{\|T(I, v, v)\|}$ is a version of **gradient ascent**.

Local Optima: $w^\top \nabla^2 L(v, \lambda) w < 0$ for all $w \perp v$, at a stationary point v

- v_1 is the only local optimum.
- All other eigenvectors are saddle points.
- $\{v_h\}_{h=1}^K$ are the only local optima.
- All spurious eigenvectors are saddle points.

Question: What about performance under noise?

Tensor Perturbation Analysis

$$\hat{\mathcal{T}} = \mathcal{T} + \mathcal{E}, \quad \mathcal{T} = \sum_h \lambda_h \mathbf{v}_h \otimes^3, \quad \|\mathcal{E}\| := \max_{\mathbf{x}: \|\mathbf{x}\|=1} |\mathcal{E}(\mathbf{x}, \mathbf{x}, \mathbf{x})| \leq \epsilon.$$

Tensor Perturbation Analysis

$$\hat{\mathcal{T}} = \mathcal{T} + \mathcal{E}, \quad \mathcal{T} = \sum_h \lambda_h \mathbf{v}_h \otimes^3, \quad \|\mathcal{E}\| := \max_{\mathbf{x}: \|\mathbf{x}\|=1} |\mathcal{E}(\mathbf{x}, \mathbf{x}, \mathbf{x})| \leq \epsilon.$$

Theorem: Let T be number of iterations. If

$$T \geq \log K + \log \log \frac{\lambda_{\max}}{\epsilon}, \quad \epsilon < \frac{\lambda_{\min}}{K},$$

then output (\mathbf{v}, λ) (after **polynomial restarts**) satisfies

$$\|\mathbf{v} - \mathbf{v}_1\| \leq O\left(\frac{\epsilon}{\lambda_1}\right), \quad \|\lambda - \lambda_1\| \leq O(\epsilon),$$

where \mathbf{v}_1 is s.t. $\lambda_1 |c_1| > \lambda_2 |c_2| \dots$, $c_i := \langle \mathbf{v}_i, \mathbf{u}_0 \rangle$, and \mathbf{u}_0 is the (successful) initializer.

Tensor Perturbation Analysis

$$\hat{\mathcal{T}} = \mathcal{T} + \mathcal{E}, \quad \mathcal{T} = \sum_h \lambda_h \mathbf{v}_h \otimes^3, \quad \|\mathcal{E}\| := \max_{\mathbf{x}: \|\mathbf{x}\|=1} |\mathcal{E}(\mathbf{x}, \mathbf{x}, \mathbf{x})| \leq \epsilon.$$

Theorem: Let T be number of iterations. If

$$T \geq \log K + \log \log \frac{\lambda_{\max}}{\epsilon}, \quad \epsilon < \frac{\lambda_{\min}}{K},$$

then output (\mathbf{v}, λ) (after **polynomial restarts**) satisfies

$$\|\mathbf{v} - \mathbf{v}_1\| \leq O\left(\frac{\epsilon}{\lambda_1}\right), \quad \|\lambda - \lambda_1\| \leq O(\epsilon),$$

where \mathbf{v}_1 is s.t. $\lambda_1 |c_1| > \lambda_2 |c_2| \dots$, $c_i := \langle \mathbf{v}_i, \mathbf{u}_0 \rangle$, and \mathbf{u}_0 is the (successful) initializer.

- Careful analysis of deflation: avoid buildup of errors.
- Implies **polynomial sample complexity** for learning.

Other tensor decomposition techniques

Orthogonal Tensor Decomposition

Simultaneous Power Method

- (Wang & Lu, 2017)
 - Simultaneous recovery of eigenvectors
 - Initialization is not optimal

Orthogonalized Simultaneous Alternating Least Square

- (Sharan & Valiant, 2017)
 - Random initialization
 - Proved convergence for symmetric tensor

Initialization

- SVD based initialization (Anandkumar & Janzamin, 2014).
- State-of-the-art (trace based) initialization (Li & Huang, 2018).

Outline

- 1 Introduction
- 2 Motivation: Challenges of MLE for Gaussian Mixtures
- 3 Introduction of Method of Moments and Tensor Notations
- 4 Topic Model for Single-topic Documents
- 5 Algorithms for Tensor Decompositions
- 6 Tensor Decomposition for Neural Network Compression**
- 7 Conclusion

Neural Network - Nonlinear Function Approximation



Image classification

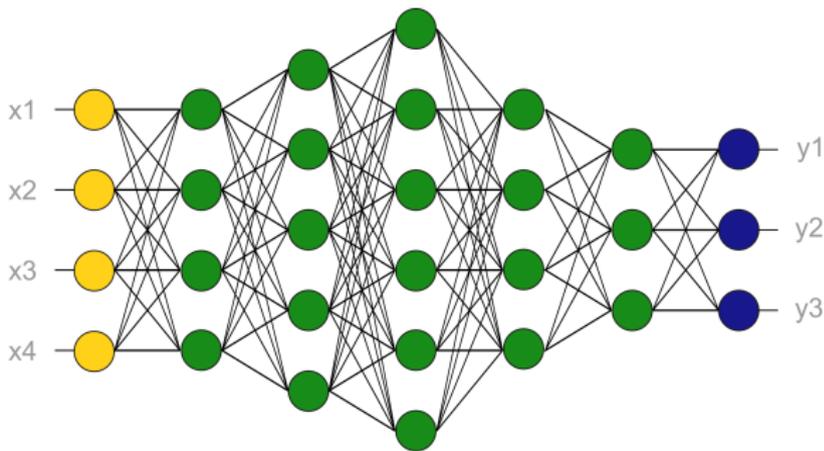


Speech recognition



Text processing

Success of Deep Neural Networks



- computation power growth
- enormous labeled data

Neural Network - Nonlinear Function Approximation



Image classification

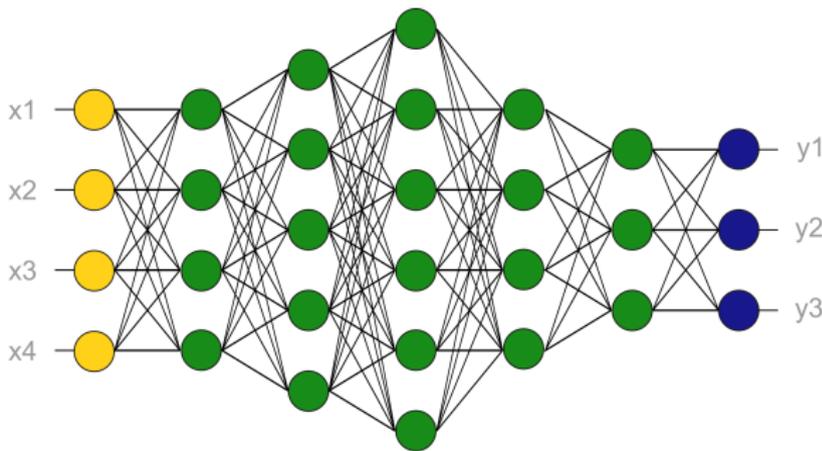


Speech recognition



Text processing

Success of Deep Neural Networks

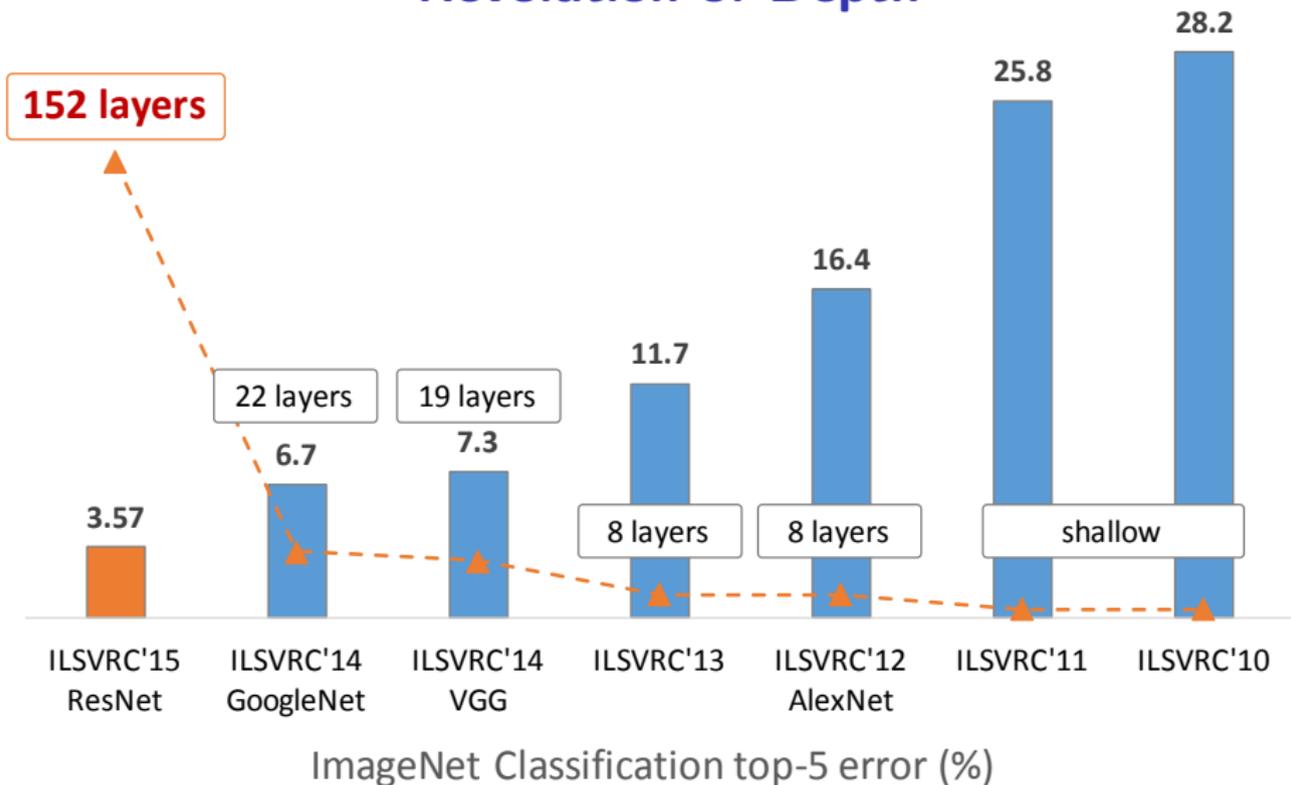


- computation power growth
- enormous labeled data

Express Power

- linear composition vs nonlinear composition
- shallow network vs deep structure

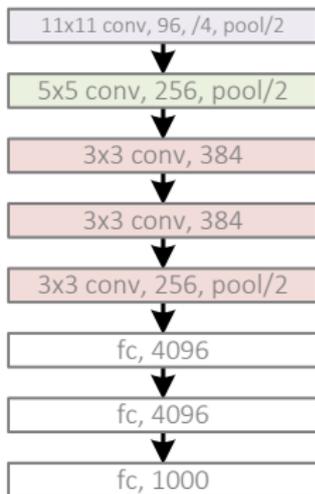
Revolution of Depth



Kaiming He, Xiangyu Zhang, Shaoqing Ren, & Jian Sun. "Deep Residual Learning for Image Recognition". CVPR 2016.

Revolution of Depth

AlexNet, 8 layers
(ILSVRC 2012)



Kaiming He, Xiangyu Zhang, Shaoqing Ren, & Jian Sun. "Deep Residual Learning for Image Recognition". CVPR 2016.

Revolution of Depth

AlexNet, 8 layers
(ILSVRC 2012)



VGG, 19 layers
(ILSVRC 2014)

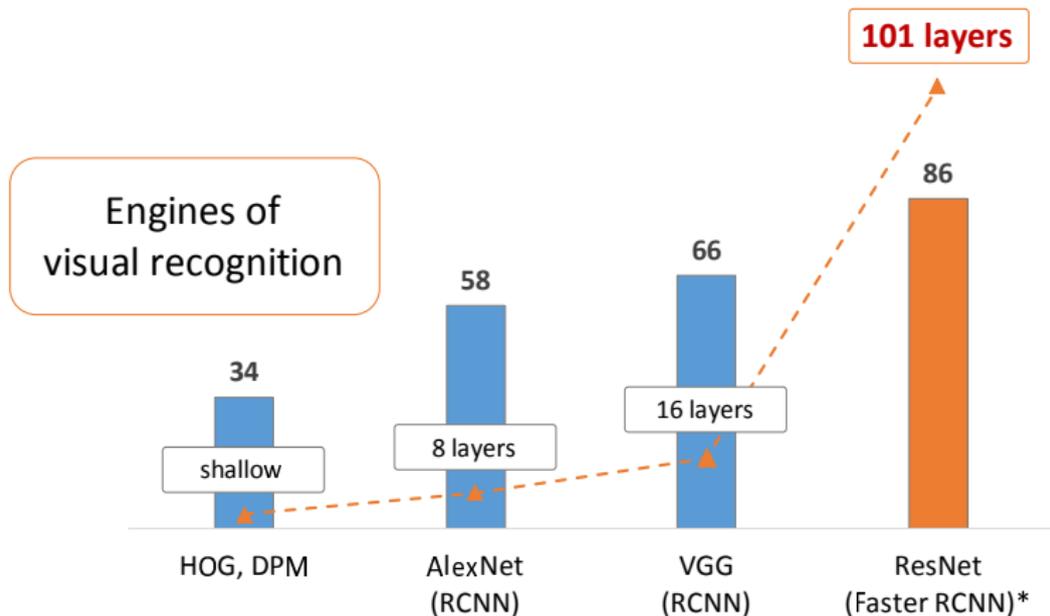


ResNet, **152 layers**
(ILSVRC 2015)



Kaiming He, Xiangyu Zhang, Shaoqing Ren, & Jian Sun. "Deep Residual Learning for Image Recognition". CVPR 2016.

Revolution of Depth



PASCAL VOC 2007 **Object Detection** mAP (%)

*w/ other improvements & more data

Kaiming He, Xiangyu Zhang, Shaoqing Ren, & Jian Sun. "Deep Residual Learning for Image Recognition". CVPR 2016.

Challenges For Large Deep Neural Network

Learning

- Learning takes longer, might not converge, susceptible to vanishing/exploding gradients, etc
- One-time cost.

Challenges For Large Deep Neural Network

Learning

- Learning takes longer, might not converge, susceptible to vanishing/exploding gradients, etc
- One-time cost.

Test

- Requires large amount of computation and memory storage.
Ill-suited for smart phones or IoT device.
- Repeated cost.

Challenges For Large Deep Neural Network

Learning

- Learning takes longer, might not converge, susceptible to vanishing/exploding gradients, etc
- One-time cost.

Test

- Requires large amount of computation and memory storage.
Ill-suited for smart phones or IoT device.
- Repeated cost.

How to compress the neural network without much performance loss?

Common Types of Tensor Decompositions

m -order tensor $\mathcal{T} \in \mathbb{R}^{I_0 \times I_1 \times \cdots \times I_{m-1}}$

Common Types of Tensor Decompositions

m -order tensor $\mathcal{T} \in \mathbb{R}^{I_0 \times I_1 \times \dots \times I_{m-1}}$

CANDECOMP/PARAFAC (CP) Decomposition

- Factorize a tensor into sum of rank-1 tensors
- Rank-1 tensor is defined as outer product of multiple vectors
- $\mathcal{T}_{i_0, \dots, i_{m-1}} = \sum_{r=0}^{R-1} M_{r, i_0}^{(0)} \cdots M_{r, i_{m-1}}^{(m-1)}$

Common Types of Tensor Decompositions

m -order tensor $\mathcal{T} \in \mathbb{R}^{I_0 \times I_1 \times \dots \times I_{m-1}}$

CANDECOMP/PARAFAC (CP) Decomposition

- Factorize a tensor into sum of rank-1 tensors
- Rank-1 tensor is defined as outer product of multiple vectors
- $\mathcal{T}_{i_0, \dots, i_{m-1}} = \sum_{r=0}^{R-1} M_{r, i_0}^{(0)} \dots M_{r, i_{m-1}}^{(m-1)}$

Tucker (TK) Decomposition

- More general than CP decomposition
- Multilinear operation on a core tensor \mathcal{C} : $\mathcal{C}(M^{(0)}, \dots, M^{(m-1)})$
- $\mathcal{T}_{i_0, \dots, i_{m-1}} = \sum_{r_0=0}^{R_0-1} \dots \sum_{r_{m-1}=0}^{R_{m-1}-1} \mathcal{C}_{r_0, \dots, r_{m-1}} M_{r_0, i_0}^{(0)} \dots M_{r_{m-1}, i_{m-1}}^{(m-1)}$

Common Types of Tensor Decompositions

m -order tensor $\mathcal{T} \in \mathbb{R}^{I_0 \times I_1 \times \dots \times I_{m-1}}$

CANDECOMP/PARAFAC (CP) Decomposition

- Factorize a tensor into sum of rank-1 tensors
- Rank-1 tensor is defined as outer product of multiple vectors
- $\mathcal{T}_{i_0, \dots, i_{m-1}} = \sum_{r=0}^{R-1} M_{r, i_0}^{(0)} \dots M_{r, i_{m-1}}^{(m-1)}$

Tucker (TK) Decomposition

- More general than CP decomposition
- Multilinear operation on a core tensor \mathcal{C} : $\mathcal{C}(M^{(0)}, \dots, M^{(m-1)})$
- $\mathcal{T}_{i_0, \dots, i_{m-1}} = \sum_{r_0=0}^{R_0-1} \dots \sum_{r_{m-1}=0}^{R_{m-1}-1} \mathcal{C}_{r_0, \dots, r_{m-1}} M_{r_0, i_0}^{(0)} \dots M_{r_{m-1}, i_{m-1}}^{(m-1)}$

Tensor-Train (TT) Decomposition

- Factorize a tensor into a number of interconnected lower-order tensors
- $\mathcal{T}_{i_0, \dots, i_{m-1}} = \sum_{r_0=1}^{R_0-1} \dots \sum_{r_{m-2}=1}^{R_{m-2}-1} \mathcal{T}_{i_0, r_0}^{(0)} \mathcal{T}_{r_0, i_1, r_1}^{(1)} \dots \mathcal{T}_{r_{m-2}, i_{m-1}}^{(m-1)}$

Compression of Convolutional Layer w/ Tensor Decompositions

Convolutional Kernel: tensor $\mathcal{K} \in \mathbb{R}^{H \times W \times S \times T}$

Compression of Convolutional Layer w/ Tensor Decompositions

Convolutional Kernel: tensor $\mathcal{K} \in \mathbb{R}^{H \times W \times S \times T}$

- Filter height/width H/W , No. of input/output channels S/T .

Compression of Convolutional Layer w/ Tensor Decompositions

Convolutional Kernel: tensor $\mathcal{K} \in \mathbb{R}^{H \times W \times S \times T}$

- Filter height/width H/W , No. of input/output channels S/T .
- Map an input tensor $\mathcal{U} \in \mathbb{R}^{X \times Y \times S}$ to an output tensor $\mathcal{V} \in \mathbb{R}^{X' \times Y' \times T}$.

Compression of Convolutional Layer w/ Tensor Decompositions

Convolutional Kernel: tensor $\mathcal{K} \in \mathbb{R}^{H \times W \times S \times T}$

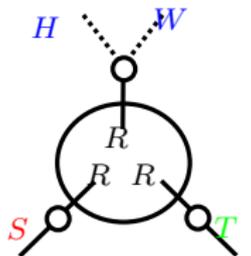
- Filter height/width H/W , No. of input/output channels S/T .
- Map an input tensor $\mathcal{U} \in \mathbb{R}^{X \times Y \times S}$ to an output tensor $\mathcal{V} \in \mathbb{R}^{X' \times Y' \times T}$.

Kernel CP Decomposition

- **CP:** Decompose kernel \mathcal{K} into 3 factor tensors

$$\mathcal{K}_{i,j,s,t} = \sum_{r=0}^{R-1} \mathcal{K}_{s,r}^{(0)} \mathcal{K}_{i,j,r}^{(1)} \mathcal{K}_{r,t}^{(2)}$$

- No. of param.: $HWST \rightarrow (HW + S + T)R$



CP decomposition

Compression of Convolutional Layer w/ Tensor Decompositions

Convolutional Kernel: tensor $\mathcal{K} \in \mathbb{R}^{H \times W \times S \times T}$

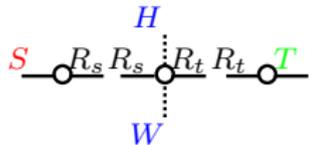
- Filter height/width H/W , No. of input/output channels S/T .
- Map an input tensor $\mathcal{U} \in \mathbb{R}^{X \times Y \times S}$ to an output tensor $\mathcal{V} \in \mathbb{R}^{X' \times Y' \times T}$.

Kernel TK Decomposition

- **TK:** Decompose \mathcal{K} into 1 core tensor, 2 factor tensors

$$\mathcal{K}_{i,j,s,t} = \sum_{r_s=0}^{R_s-1} \sum_{r_t=0}^{R_t-1} \mathcal{K}_{s,r_s}^{(0)} \mathcal{K}_{i,j,r_s,r_t}^{(1)} \mathcal{K}_{r_t,t}^{(2)}$$

- No. of param.: $HWST \rightarrow SR_s + HWR_sR_t + R_tT$



TK decomposition

Compression of Convolutional Layer w/ Tensor Decompositions

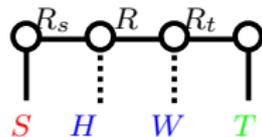
Convolutional Kernel: tensor $\mathcal{K} \in \mathbb{R}^{H \times W \times S \times T}$

- Filter height/width H/W , No. of input/output channels S/T .
- Map an input tensor $\mathcal{U} \in \mathbb{R}^{X \times Y \times S}$ to an output tensor $\mathcal{V} \in \mathbb{R}^{X' \times Y' \times T}$.

Kernel TT Decomposition

- **TT:** Decompose \mathcal{K} into 4 factor tensors

$$\mathcal{K}_{i,j,s,t} = \sum_{r_s=0}^{R_s-1} \sum_{r=0}^{R-1} \sum_{r_t=0}^{R_t-1} \mathcal{K}_{s,r_s}^{(0)} \mathcal{K}_{r_s,i,r}^{(1)} \mathcal{K}_{r,j,r_t}^{(2)} \mathcal{K}_{r_t,t}^{(3)}$$



TT decomposition

- No. of param.: $HWST \rightarrow SR_s + HR_sR + WR_tR + R_tT$

Tensorized Spectrum Preserving Compression of Neural Networks

Convolutional Kernel: $\mathcal{K} \in \mathbb{R}^{H \times W \times S \times T}$ tensorized to
 $\mathcal{K}' \in \mathbb{R}^{H \times W \times S_0 \times \dots \times S_{m-1} \times T_0 \times \dots \times T_{m-1}}$

Tensorized Spectrum Preserving Compression of Neural Networks

Convolutional Kernel: $\mathcal{K} \in \mathbb{R}^{H \times W \times S \times T}$ **tensorized** to
 $\mathcal{K}' \in \mathbb{R}^{H \times W \times S_0 \times \dots \times S_{m-1} \times T_0 \times \dots \times T_{m-1}}$

- Tensorization: kernel reshaped to higher order tensor.

Tensorized Spectrum Preserving Compression of Neural Networks

Convolutional Kernel: $\mathcal{K} \in \mathbb{R}^{H \times W \times S \times T}$ **tensorized** to
 $\mathcal{K}' \in \mathbb{R}^{H \times W \times S_0 \times \dots \times S_{m-1} \times T_0 \times \dots \times T_{m-1}}$

- Tensorization: kernel reshaped to higher order tensor.
- $S = \prod_{i=0}^{m-1} S_i$ and $T = \prod_{i=0}^{m-1} T_i$.

Tensorized Spectrum Preserving Compression of Neural Networks

Convolutional Kernel: $\mathcal{K} \in \mathbb{R}^{H \times W \times S \times T}$ tensorized to
 $\mathcal{K}' \in \mathbb{R}^{H \times W \times S_0 \times \dots \times S_{m-1} \times T_0 \times \dots \times T_{m-1}}$

- Tensorization: kernel reshaped to higher order tensor.
- $S = \prod_{i=0}^{m-1} S_i$ and $T = \prod_{i=0}^{m-1} T_i$.
- Input tensor $\mathcal{U} \in \mathbb{R}^{X \times Y \times S}$ tensorized to $\mathcal{U}' \in \mathbb{R}^{X \times Y \times S_0 \times \dots \times S_{m-1}}$.

Tensorized Spectrum Preserving Compression of Neural Networks

Convolutional Kernel: $\mathcal{K} \in \mathbb{R}^{H \times W \times S \times T}$ tensorized to
 $\mathcal{K}' \in \mathbb{R}^{H \times W \times S_0 \times \dots \times S_{m-1} \times T_0 \times \dots \times T_{m-1}}$

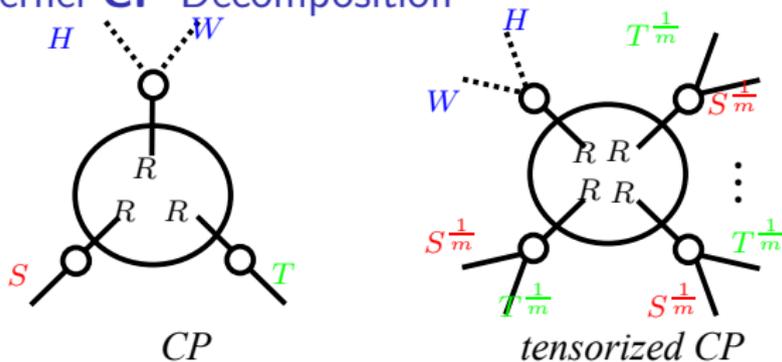
- Tensorization: kernel reshaped to higher order tensor.
- $S = \prod_{i=0}^{m-1} S_i$ and $T = \prod_{i=0}^{m-1} T_i$.
- Input tensor $\mathcal{U} \in \mathbb{R}^{X \times Y \times S}$ tensorized to $\mathcal{U}' \in \mathbb{R}^{X \times Y \times S_0 \times \dots \times S_{m-1}}$.
- Output reshaped $\mathcal{V} \in \mathbb{R}^{X \times Y \times T}$ to $\mathcal{V}' \in \mathbb{R}^{X' \times Y' \times T_0 \times \dots \times T_{m-1}}$.

Tensorized Spectrum Preserving Compression of Neural Networks

Convolutional Kernel: $\mathcal{K} \in \mathbb{R}^{H \times W \times S \times T}$ **tensorized** to
 $\mathcal{K}' \in \mathbb{R}^{H \times W \times S_0 \times \dots \times S_{m-1} \times T_0 \times \dots \times T_{m-1}}$

- Tensorization: kernel reshaped to higher order tensor.
- $S = \prod_{i=0}^{m-1} S_i$ and $T = \prod_{i=0}^{m-1} T_i$.
- Input tensor $\mathcal{U} \in \mathbb{R}^{X \times Y \times S}$ tensorized to $\mathcal{U}' \in \mathbb{R}^{X \times Y \times S_0 \times \dots \times S_{m-1}}$.
- Output reshaped $\mathcal{V} \in \mathbb{R}^{X \times Y \times T}$ to $\mathcal{V}' \in \mathbb{R}^{X' \times Y' \times T_0 \times \dots \times T_{m-1}}$.

Tensorized Kernel **CP** Decomposition



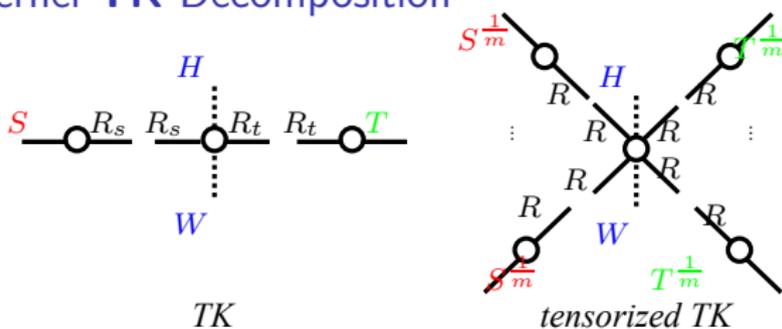
- Param. No.: $HWST \rightarrow (HW + S + T)R \rightarrow (m(ST)^{\frac{1}{m}} + HW)R$

Tensorized Spectrum Preserving Compression of Neural Networks

Convolutional Kernel: $\mathcal{K} \in \mathbb{R}^{H \times W \times S \times T}$ **tensorized** to
 $\mathcal{K}' \in \mathbb{R}^{H \times W \times S_0 \times \dots \times S_{m-1} \times T_0 \times \dots \times T_{m-1}}$

- Tensorization: kernel reshaped to higher order tensor.
- $S = \prod_{i=0}^{m-1} S_i$ and $T = \prod_{i=0}^{m-1} T_i$.
- Input tensor $\mathcal{U} \in \mathbb{R}^{X \times Y \times S}$ tensorized to $\mathcal{U}' \in \mathbb{R}^{X \times Y \times S_0 \times \dots \times S_{m-1}}$.
- Output reshaped $\mathcal{V} \in \mathbb{R}^{X \times Y \times T}$ to $\mathcal{V}' \in \mathbb{R}^{X' \times Y' \times T_0 \times \dots \times T_{m-1}}$.

Tensorized Kernel **TK** Decomposition



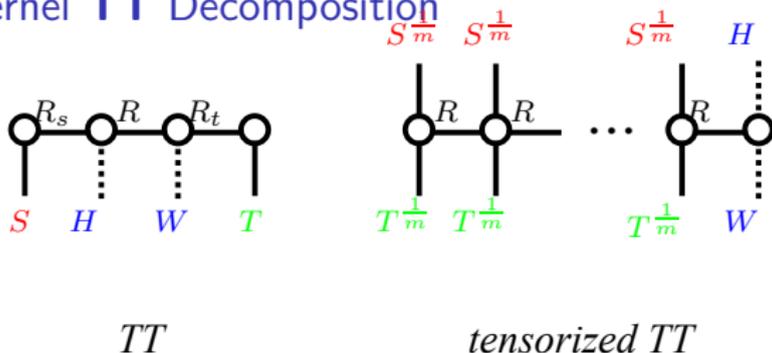
- Param. No.: $HWST \rightarrow SR_s + HW R_s R_t + R_t T \rightarrow m(S^{\frac{1}{m}} + T^{\frac{1}{m}})R + HW R^{2m}$

Tensorized Spectrum Preserving Compression of Neural Networks

Convolutional Kernel: $\mathcal{K} \in \mathbb{R}^{H \times W \times S \times T}$ tensorized to $\mathcal{K}' \in \mathbb{R}^{H \times W \times S_0 \times \dots \times S_{m-1} \times T_0 \times \dots \times T_{m-1}}$

- Tensorization: kernel reshaped to higher order tensor.
- $S = \prod_{i=0}^{m-1} S_i$ and $T = \prod_{i=0}^{m-1} T_i$.
- Input tensor $\mathcal{U} \in \mathbb{R}^{X \times Y \times S}$ tensorized to $\mathcal{U}' \in \mathbb{R}^{X \times Y \times S_0 \times \dots \times S_{m-1}}$.
- Output reshaped $\mathcal{V} \in \mathbb{R}^{X \times Y \times T}$ to $\mathcal{V}' \in \mathbb{R}^{X' \times Y' \times T_0 \times \dots \times T_{m-1}}$.

Tensorized Kernel TT Decomposition



- Param. No.: $HWST \rightarrow SR_s + HR_sR + WR_tR + R_tT \rightarrow (m(ST)^{\frac{1}{m}}R + HW)R$

Experiments - Compress CIFAR10 Resnet-34

Successful Compression of CIFAR10 Resnet-34 Network (Su, Li, Bhattacharjee & **Huang**, 2018)

Method	Compression rate: SPC, E2E				Compression rate: t-SPC , Seq.			
	5%	10%	20%	40%	2%	5%	10%	20%
CP	84.02	86.93	88.75	88.75	85.7	89.86	91.28	-
TK	83.57	86.00	88.03	89.35	61.06	71.34	81.59	87.11
TT	77.44	82.92	84.13	86.64	78.95	84.26	87.89	-

- Testing accuracies of tensor methods under compression rates.
- The uncompressed network achieves 93.2% accuracy.
- CIFAR10 Resnet-34 has 4×10^5 parameters that have to be trained and retained during testing.

Experiments - Compress ImageNet Resnet-50

Successful Compression of ImageNet Resnet-50 Network (Su, Li, Bhattacharjee & Huang, 2018)

# Epochs	Uncompressed	SPC-TT (E2E)	t-SPC-TT (Seq.)
0.2	4.22	0.66x	10.51x
0.3	6.23	0.64x	7.54x
0.5	9.01	0.83x	5.54x
1.0	17.3	0.74x	3.04x
2.0	30.8	0.59x	1.75x

- Testing accuracy of tensor methods compared to the uncompressed ImageNet Resnet-50.
- The accuracy of the tensor method results (both non-tensorized and tensorized) are shown normalized to the uncompressed network's accuracy.

Outline

- 1 Introduction
- 2 Motivation: Challenges of MLE for Gaussian Mixtures
- 3 Introduction of Method of Moments and Tensor Notations
- 4 Topic Model for Single-topic Documents
- 5 Algorithms for Tensor Decompositions
- 6 Tensor Decomposition for Neural Network Compression
- 7 Conclusion**

Conclusion

- Method-of-moments can efficiently estimate parameters for many latent variable models.
 - ▶ Exploit distributional properties, multi-view structure, and other structure to determine usable moments tensors.
 - ▶ Some efficient algorithms for carrying out the tensor decomposition to obtain parameter estimates.
- Tensor decomposition of neural network kernels/weights effectively compresses the network.
- Many issues to resolve
 - ▶ Handle model misspecification, increase robustness.
 - ▶ Learning deep neural network parameters using tensor decomposition?

A Short List of Related Papers to Today's Talk

- "A Method of Moments for Mixture Models and Hidden Markov Models", by Anima Anandkumar, Daniel Hsu and Sham Kakade. In Conference on Learning Theory, 2012.
- "Tensor Decompositions for Learning Latent Variable Models", by Anima Anandkumar, Rong Ge, Daniel Hsu, Sham Kakade and Matus Telgarsky. In Journal of Machine Learning Research, 2014.
- "Escaping from saddle points online stochastic gradient for tensor decomposition", Rong Ge, Furong Huang, Chi Jin and Yang Yuan. In Conference on Learning Theory, 2015.
- "Online tensor methods for learning latent variable models", Furong Huang, Niranjan U. N., Mohammad Umar Hakeem and Anima Anandkumar. The Journal of Machine Learning Research, 2016.
- "Guaranteed Simultaneous Asymmetric Tensor Decomposition via Orthogonalized Alternating Least Squares", by Jialin Li and Furong Huang, 2018.
- "Tensorized Spectrum Preserving Compression for Neural Networks", by Jiahao Su, Jingling Li, Bobby Bhattacharjee and Furong Huang, 2018.

Tensor Softwares

- Spark implementation of method of moments to learn Latent Dirichlet Allocation available at <https://github.com/FurongHuang/spectrallda-tensorspark>.
- Tensorly: Simple and Fast Tensor Learning in Python available at <http://tensorly.org/stable/home.html>.
- A general library with higher order tensor operations is coming soon.