

# Learning latent variable models using tensor decompositions

Daniel Hsu

Computer Science Department & Data Science Institute  
Columbia University

Machine Learning Summer School  
June 29-30, 2018

# El tema (subject matter)

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

# El tema (subject matter)

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

“Method-of-moments” (Pearson, 1894)

## Example #1: summarizing a corpus of documents

Observation: **documents express one or more thematic topics.**

### *Politics Ensnare Mohamed Salah and Switzerland at the World Cup*

By Rory Smith, James Montague and Tariq Panja

June 24, 2018



MOSCOW — The World Cup was thrust into the combustible mix of politics and soccer — dangerous ground that world soccer takes great pains to avoid — as a growing number of disciplinary proceedings and a star player's threatened retirement brought several sensitive international flash points to the tournament's doorstep this weekend.

## Example #1: summarizing a corpus of documents

Observation: **documents express one or more thematic topics.**

### *Politics Ensnare Mohamed Salah and Switzerland at the World Cup*

By Rory Smith, James Montague and Tariq Panja

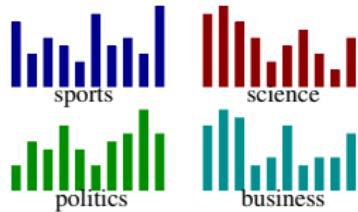
June 24, 2018



MOSCOW — The World Cup was thrust into the combustible mix of politics and soccer — dangerous ground that world soccer takes great pains to avoid — as a growing number of disciplinary proceedings and a star player's threatened retirement brought several sensitive international flash points to the tournament's doorstep this weekend.

- ▶ What topics are expressed in a corpus of documents?
- ▶ How prevalent is each topic in the corpus?

## Topic model (e.g., latent Dirichlet allocation)

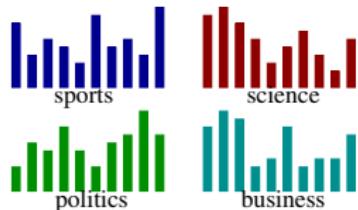


$K$  topics (distributions over vocab words).

Document  $\equiv$  mixture of topics.

Word tokens in doc.  $\stackrel{\text{iid}}{\sim}$  mixture distribution.

## Topic model (e.g., latent Dirichlet allocation)



$K$  topics (distributions over vocab words).

Document  $\equiv$  mixture of topics.

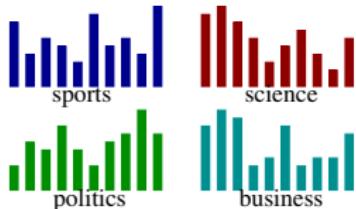
Word tokens in doc.  $\stackrel{\text{iid}}{\sim}$  mixture distribution.



E.g.,

$$\stackrel{\text{iid}}{\sim} 0.7 \times P_{\text{sports}} + 0.3 \times P_{\text{politics}}.$$

## Topic model (e.g., latent Dirichlet allocation)



$K$  topics (distributions over vocab words).

Document  $\equiv$  mixture of topics.

Word tokens in doc.  $\stackrel{\text{iid}}{\sim}$  mixture distribution.



E.g.,

$$\stackrel{\text{iid}}{\sim} 0.7 \times \mathbf{P}_{\text{sports}} + 0.3 \times \mathbf{P}_{\text{politics}}.$$

Given corpus of documents (and “hyper-parameters”, e.g.,  $K$ ), produce estimates of **model parameters**, e.g.:

- ▶ Distribution  $\mathbf{P}_t$  over vocab words, for each  $t \in [K]$ .
- ▶ Weight  $w_t$  of topic  $t$  in document corpus, for each  $t \in [K]$ .

## Labels / annotations

- ▶ Suppose each word token  $x$  in document is *annotated* with source topic  $t_x \in \{1, 2, \dots, K\}$ .

Politics	Ensnare	Mohamed_Salah	and	Switzerland	at
3	3	1	5	3	5

## Labels / annotations

- ▶ Suppose each word token  $x$  in document is *annotated* with source topic  $t_x \in \{1, 2, \dots, K\}$ .

Politics	Ensnare	Mohamed_Salah	and	Switzerland	at
3	3	1	5	3	5

Then estimating the  $\{(\mathbf{P}_t, w_t)\}_{t=1}^K$  can be done “directly”.

## Labels / annotations

- ▶ Suppose each word token  $x$  in document is *annotated* with source topic  $t_x \in \{1, 2, \dots, K\}$ .

Politics	Ensnare	Mohamed_Salah	and	Switzerland	at
3	3	1	5	3	5

Then estimating the  $\{(\mathbf{P}_t, w_t)\}_{t=1}^K$  can be done “directly”.

- ▶ Unfortunately, we often don't have such annotations  
(i.e., data are *unlabeled* / topics are *hidden*).

“Direct” approach to estimation unavailable.

## Example #2: subpopulations in data



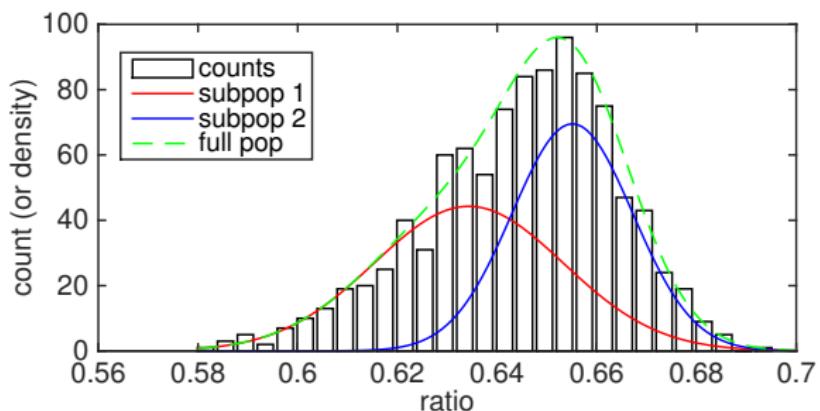
**Data studied by Pearson (1894):**  
ratio of forehead-width to body-length for 1000 crabs.

## Example #2: subpopulations in data



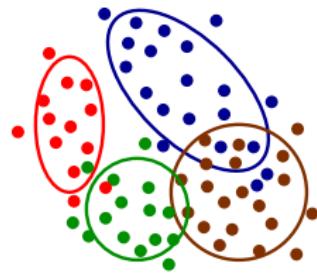
**Data studied by Pearson (1894):**  
ratio of forehead-width to body-length for 1000 crabs.

Sample may be comprised of different sub-species of crabs.



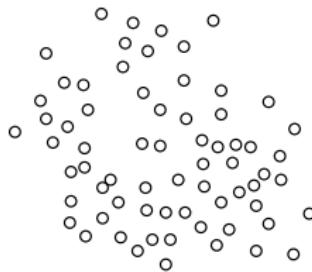
# Gaussian mixture model

$$H \sim \text{Categorical}(\pi_1, \pi_2, \dots, \pi_K);$$
$$\mathbf{X} | H = t \sim \text{Normal}(\boldsymbol{\mu}_t, \boldsymbol{\Sigma}_t), \quad t \in [K].$$



## Gaussian mixture model

$$H \sim \text{Categorical}(\pi_1, \pi_2, \dots, \pi_K);$$
$$\mathbf{X} | H = t \sim \text{Normal}(\boldsymbol{\mu}_t, \boldsymbol{\Sigma}_t), \quad t \in [K].$$



Estimate **mean vector**, **covariance matrix**, and **mixing weight** of each subpopulation from *unlabeled data*.

## Maximum likelihood estimation

- ▶ No “direct” estimators when some variables are hidden.

## Maximum likelihood estimation

- ▶ No “direct” estimators when some variables are hidden.
- ▶ **Maximum likelihood estimator (MLE):**

$$\theta_{\text{MLE}} := \arg \max_{\theta \in \Theta} \log \Pr_{\theta}(\text{data}) .$$

## Maximum likelihood estimation

- ▶ No “direct” estimators when some variables are hidden.
- ▶ **Maximum likelihood estimator (MLE):**

$$\theta_{\text{MLE}} := \arg \max_{\theta \in \Theta} \log \Pr_{\theta}(\text{data}) .$$

- ▶ **Note:** log-likelihood is not necessarily concave function of  $\theta$ .

# Maximum likelihood estimation

- ▶ No “direct” estimators when some variables are hidden.
- ▶ **Maximum likelihood estimator (MLE):**

$$\theta_{\text{MLE}} := \arg \max_{\theta \in \Theta} \log \Pr_{\theta}(\text{data}) .$$

- ▶ **Note:** log-likelihood is not necessarily concave function of  $\theta$ .
- ▶ For latent variable models, often use local optimization, most notably via **Expectation-Maximization (EM)** (Dempster, Laird, & Rubin, 1977).

## MLE for Gaussian mixture models

Given data  $\{\mathbf{x}_i\}_{i=1}^n$ , find  $\{(\boldsymbol{\mu}_t, \boldsymbol{\Sigma}_t, \pi_t)\}_{t=1}^K$  to maximize

$$\sum_{i=1}^n \log \left( \sum_{t=1}^K \pi_t \cdot \frac{1}{\det(\boldsymbol{\Sigma}_t)^{1/2}} \exp \left\{ -\frac{1}{2} (\mathbf{x}_i - \boldsymbol{\mu}_t)^\top \boldsymbol{\Sigma}_t^{-1} (\mathbf{x}_i - \boldsymbol{\mu}_t) \right\} \right).$$

## MLE for Gaussian mixture models

Given data  $\{\mathbf{x}_i\}_{i=1}^n$ , find  $\{(\boldsymbol{\mu}_t, \boldsymbol{\Sigma}_t, \pi_t)\}_{t=1}^K$  to maximize

$$\sum_{i=1}^n \log \left( \sum_{t=1}^K \pi_t \cdot \frac{1}{\det(\boldsymbol{\Sigma}_t)^{1/2}} \exp \left\{ -\frac{1}{2} (\mathbf{x}_i - \boldsymbol{\mu}_t)^\top \boldsymbol{\Sigma}_t^{-1} (\mathbf{x}_i - \boldsymbol{\mu}_t) \right\} \right).$$

- ▶ Sensible with restrictions on  $\boldsymbol{\Sigma}_t$  (e.g.,  $\boldsymbol{\Sigma}_t \succeq \sigma^2 \mathbf{I}$ ).

## MLE for Gaussian mixture models

Given data  $\{\mathbf{x}_i\}_{i=1}^n$ , find  $\{(\boldsymbol{\mu}_t, \boldsymbol{\Sigma}_t, \pi_t)\}_{t=1}^K$  to maximize

$$\sum_{i=1}^n \log \left( \sum_{t=1}^K \pi_t \cdot \frac{1}{\det(\boldsymbol{\Sigma}_t)^{1/2}} \exp \left\{ -\frac{1}{2} (\mathbf{x}_i - \boldsymbol{\mu}_t)^\top \boldsymbol{\Sigma}_t^{-1} (\mathbf{x}_i - \boldsymbol{\mu}_t) \right\} \right).$$

- ▶ Sensible with restrictions on  $\boldsymbol{\Sigma}_t$  (e.g.,  $\boldsymbol{\Sigma}_t \succeq \sigma^2 \mathbf{I}$ ).
- ▶ But **NP-hard** to maximize (Tosh and Dasgupta, 2018):  
Can't expect efficient algorithms to work for all data sets.

## Parameter learning objective

Suppose iid sample of size  $n$  is generated by distribution from model with (unknown) parameters  $\theta \in \Theta \subseteq \mathbb{R}^p$ . ( $p = \# \text{ params}$ )

## Parameter learning objective

Suppose iid sample of size  $n$  is generated by distribution from model with (unknown) parameters  $\boldsymbol{\theta} \in \Theta \subseteq \mathbb{R}^p$ . ( $p = \# \text{ params}$ )

**Task:** Produce estimate  $\hat{\boldsymbol{\theta}}$  of  $\boldsymbol{\theta}$  such that

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

(i.e.,  $\hat{\boldsymbol{\theta}}$  is *consistent*).

## Parameter learning objective

Suppose iid sample of size  $n$  is generated by distribution from model with (unknown) parameters  $\boldsymbol{\theta} \in \Theta \subseteq \mathbb{R}^p$ . ( $p = \# \text{ params}$ )

**Task:** Produce estimate  $\hat{\boldsymbol{\theta}}$  of  $\boldsymbol{\theta}$  such that

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

(i.e.,  $\hat{\boldsymbol{\theta}}$  is *consistent*).

- ▶ E.g., for spherical Gaussian mixtures:
  - ▶ For  $K = 2$  (and  $\pi_t = 1/2$ ,  $\boldsymbol{\Sigma}_t = \mathbf{I}$ ): EM is consistent (Xu, H., & Maleki, 2016).

## Parameter learning objective

Suppose iid sample of size  $n$  is generated by distribution from model with (unknown) parameters  $\boldsymbol{\theta} \in \Theta \subseteq \mathbb{R}^p$ . ( $p = \# \text{ params}$ )

**Task:** Produce estimate  $\hat{\boldsymbol{\theta}}$  of  $\boldsymbol{\theta}$  such that

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

(i.e.,  $\hat{\boldsymbol{\theta}}$  is *consistent*).

- ▶ E.g., for spherical Gaussian mixtures:
  - ▶ For  $K = 2$  (and  $\pi_t = 1/2$ ,  $\boldsymbol{\Sigma}_t = \mathbf{I}$ ): EM is consistent (Xu, H., & Maleki, 2016).
  - ▶ Larger  $K$ : easily trapped in local maxima, far from global max (Jin, Zhang, Balakrishnan, Wainwright, & Jordan, 2016).

## Parameter learning objective

Suppose iid sample of size  $n$  is generated by distribution from model with (unknown) parameters  $\boldsymbol{\theta} \in \Theta \subseteq \mathbb{R}^p$ . ( $p = \# \text{ params}$ )

**Task:** Produce estimate  $\hat{\boldsymbol{\theta}}$  of  $\boldsymbol{\theta}$  such that

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

(i.e.,  $\hat{\boldsymbol{\theta}}$  is *consistent*).

- ▶ E.g., for spherical Gaussian mixtures:
  - ▶ For  $K = 2$  (and  $\pi_t = 1/2$ ,  $\boldsymbol{\Sigma}_t = \mathbf{I}$ ): EM is consistent (Xu, H., & Maleki, 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.

## Parameter learning objective

Suppose iid sample of size  $n$  is generated by distribution from model with (unknown) parameters  $\boldsymbol{\theta} \in \Theta \subseteq \mathbb{R}^p$ . ( $p = \# \text{ params}$ )

**Task:** Produce estimate  $\hat{\boldsymbol{\theta}}$  of  $\boldsymbol{\theta}$  such that

$$\Pr\left(\|\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}\| \leq \epsilon\right) \geq 1 - \delta$$

with  $\text{poly}(p, 1/\epsilon, 1/\delta, \dots)$  sample size and running time.

- ▶ E.g., for spherical Gaussian mixtures:
  - ▶ For  $K = 2$  (and  $\pi_t = 1/2$ ,  $\boldsymbol{\Sigma}_t = \mathbf{I}$ ): EM is consistent (Xu, H., & Maleki, 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.

# Barriers

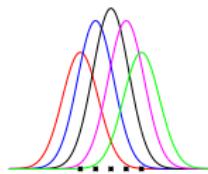
Hard to learn model parameters,  
even when data is generated by a model distribution.

# Barriers

Hard to learn model parameters,  
even when data is generated by a model distribution.



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 barriers

- ▶ **Separation conditions.**

E.g., assume mixture component distributions are far apart.

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

## Ways around the barriers

- ▶ **Separation conditions.**

E.g., assume mixture component distributions are far apart.

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

- ▶ **Structural assumptions.**

E.g., sparsity, anchor words.

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

## Ways around the barriers

- ▶ **Separation conditions.**

E.g., assume mixture component distributions are far apart.

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

- ▶ **Structural assumptions.**

E.g., sparsity, anchor words.

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

- ▶ **Non-degeneracy conditions.**

E.g., assume  $\mu_1, \mu_2, \dots, \mu_K$  are in general position.

# Ways around the barriers

- ▶ **Separation conditions.**

E.g., assume mixture component distributions are far apart.

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

- ▶ **Structural assumptions.**

E.g., sparsity, anchor words.

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

- ▶ **Non-degeneracy conditions.**

E.g., assume  $\mu_1, \mu_2, \dots, \mu_K$  are in general position.

**This lecture:** learning algorithms for non-degenerate instances via *method-of-moments*.

## Method-of-moments at a glance

1. Determine function of model parameters  $\theta$  estimatable from observable data:

$$\mathbb{E}_\theta[f(\mathbf{X})] \quad (\text{"moments"}) .$$

2. Form estimates of moments using data (e.g., iid sample):

$$\widehat{\mathbb{E}}[f(\mathbf{X})] \quad (\text{"empirical moments"}) .$$

3. Approximately solve equations for parameters  $\theta$ :

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

4. ("Fine-tune" estimated parameters with local optimization.)

# Method-of-moments at a glance

1. Determine function of model parameters  $\theta$  estimatable from observable data:

$$\mathbb{E}_\theta[f(\mathbf{X})] \quad (\text{"moments"}) .$$

**Which moments?**

2. Form estimates of moments using data (e.g., iid sample):

$$\widehat{\mathbb{E}}[f(\mathbf{X})] \quad (\text{"empirical moments"}) .$$

3. Approximately solve equations for parameters  $\theta$ :

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

**How?**

4. ("Fine-tune" estimated parameters with local optimization.)

## Method-of-moments at a glance

1. Determine function of model parameters  $\theta$  estimatable from observable data:

$$\mathbb{E}_\theta[f(\mathbf{X})] \quad (\text{"moments"}) .$$

**Which moments? Often low-order moments suffice.**

2. Form estimates of moments using data (e.g., iid sample):

$$\widehat{\mathbb{E}}[f(\mathbf{X})] \quad (\text{"empirical moments"}) .$$

3. Approximately solve equations for parameters  $\theta$ :

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

**How? Algorithms for tensor decomposition.**

4. ("Fine-tune" estimated parameters with local optimization.)

## A simple example of the method-of-moments

Let  $X \sim \text{Normal}(\mu, \sigma^2)$ . How to estimate  $\sigma^2$  from iid sample?

## A simple example of the method-of-moments

Let  $X \sim \text{Normal}(\mu, \sigma^2)$ . How to estimate  $\sigma^2$  from iid sample?

- ▶ Consider first- and second-moments:  $\mathbb{E}[X]$  and  $\mathbb{E}[X^2]$ .

## A simple example of the method-of-moments

Let  $X \sim \text{Normal}(\mu, \sigma^2)$ . How to estimate  $\sigma^2$  from iid sample?

- ▶ Consider first- and second-moments:  $\mathbb{E}[X]$  and  $\mathbb{E}[X^2]$ .
- ▶ Formula for  $\sigma^2$  in terms of moments:

$$\mathbb{E}[X^2] - \mathbb{E}[X]^2 = (\sigma^2 + \mu^2) - \mu^2 = \sigma^2.$$

## A simple example of the method-of-moments

Let  $X \sim \text{Normal}(\mu, \sigma^2)$ . How to estimate  $\sigma^2$  from iid sample?

- ▶ Consider first- and second-moments:  $\mathbb{E}[X]$  and  $\mathbb{E}[X^2]$ .
- ▶ Formula for  $\sigma^2$  in terms of moments:

$$\mathbb{E}[X^2] - \mathbb{E}[X]^2 = (\sigma^2 + \mu^2) - \mu^2 = \sigma^2.$$

- ▶ Form estimates of  $\mathbb{E}[X]$  and  $\mathbb{E}[X^2]$  from iid sample  $\{x_i\}_{i=1}^n$ :  
e.g.,

$$\widehat{\mathbb{E}}[X] := \frac{1}{n} \sum_{i=1}^n x_i, \quad \widehat{\mathbb{E}}[X^2] := \frac{1}{n} \sum_{i=1}^n x_i^2.$$

## A simple example of the method-of-moments

Let  $X \sim \text{Normal}(\mu, \sigma^2)$ . How to estimate  $\sigma^2$  from iid sample?

- ▶ Consider first- and second-moments:  $\mathbb{E}[X]$  and  $\mathbb{E}[X^2]$ .
- ▶ Formula for  $\sigma^2$  in terms of moments:

$$\mathbb{E}[X^2] - \mathbb{E}[X]^2 = (\sigma^2 + \mu^2) - \mu^2 = \sigma^2.$$

- ▶ Form estimates of  $\mathbb{E}[X]$  and  $\mathbb{E}[X^2]$  from iid sample  $\{x_i\}_{i=1}^n$ :  
e.g.,

$$\hat{\mathbb{E}}[X] := \frac{1}{n} \sum_{i=1}^n x_i, \quad \hat{\mathbb{E}}[X^2] := \frac{1}{n} \sum_{i=1}^n x_i^2.$$

- ▶ Then estimate  $\sigma^2$  with

$$\hat{\sigma}^2 := \hat{\mathbb{E}}[X^2] - \hat{\mathbb{E}}[X]^2.$$

## A simple example of the method-of-moments

Let  $X \sim \text{Normal}(\mu, \sigma^2)$ . How to estimate  $\sigma^2$  from iid sample?

- ▶ Consider first- and second-moments:  $\mathbb{E}[X]$  and  $\mathbb{E}[X^2]$ .
- ▶ Formula for  $\sigma^2$  in terms of moments:

$$\mathbb{E}[X^2] - \mathbb{E}[X]^2 = (\sigma^2 + \mu^2) - \mu^2 = \sigma^2.$$

- ▶ Form estimates of  $\mathbb{E}[X]$  and  $\mathbb{E}[X^2]$  from iid sample  $\{x_i\}_{i=1}^n$ :  
e.g.,

$$\hat{\mathbb{E}}[X] := \frac{1}{n} \sum_{i=1}^n x_i, \quad \hat{\mathbb{E}}[X^2] := \frac{1}{n} \sum_{i=1}^n x_i^2.$$

- ▶ Then estimate  $\sigma^2$  with

$$\hat{\sigma}^2 := \hat{\mathbb{E}}[X^2] - \hat{\mathbb{E}}[X]^2.$$

We'll follow this same basic recipe for much richer models!

# Outline

1. Topic model for single-topic documents.
  - ▶ Identifiability.
  - ▶ Parameter recovery via orthogonal tensor decomposition.
2. Moment decompositions for other models.
  - ▶ Mixtures of Gaussians and linear regressions.
  - ▶ Multi-view models (e.g., HMMs).
  - ▶ Other models (e.g., single-index models).
3. Error analysis (if time).

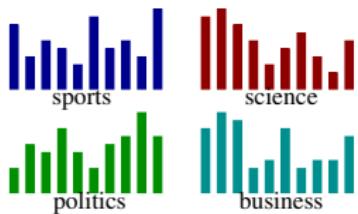
## Other models amenable to moment tensor decomposition

- ▶ Models for independent components analysis (Comon, 1994; Frieze, Jerrum, & Kannan, 1996; Arora, Ge, Moitra & Sachdeva, 2012; Anandkumar, Foster, H., Kakade, & Liu, 2012, 2015; Belkin, Rademacher, & Voss, 2013; etc.)
- ▶ Latent Dirichlet Allocation (Anandkumar, Foster, H., Kakade, & Liu, 2012, 2015; Anderson, Goyal, & Rademacher, 2013)
- ▶ Mixed-membership stochastic blockmodels (Anandkumar, Ge, H., & Kakade, 2013, 2014)
- ▶ Simple probabilistic grammars (H., Kakade, & Liang, 2012)
- ▶ Noisy-or networks (Halpern & Sontag, 2013; Jernite, Halpern & Sontag, 2013; Arora, Ge, Ma, & Risteski, 2016)
- ▶ Indian buffet process (Tung & Smola, 2014)
- ▶ Mixed multinomial logit model (Oh & Shah, 2014)
- ▶ Dawid-Skene model (Zhang, Chen, Zhou, & Jordan, 2014)
- ▶ Multi-task bandits (Azar, Lazaric, & Brunskill, 2013)
- ▶ Partially obs. MDPs (Azizzadenesheli, Lazaric, & Anandkumar, 2016)
- ▶ ...

## 1. Topic model for single-topic documents

# Topic model

## General topic model (e.g., Latent Dirichlet Allocation)



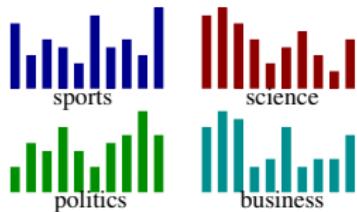
$K$  topics (dists. over words)  $\{\mathbf{P}_t\}_{t=1}^K$ .

Document  $\equiv$  mixture of topics (**hidden**).

Word tokens in doc.  $\stackrel{\text{iid}}{\sim}$  mixture distribution.

# Topic model

## Topic model for single-topic documents



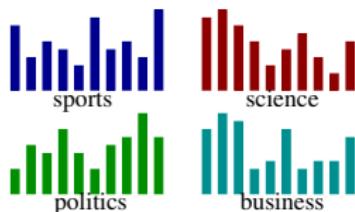
$K$  topics (dists. over words)  $\{\mathbf{P}_t\}_{t=1}^K$ .

Pick topic  $t$  with prob.  $w_t$  (hidden).

Word tokens in doc.  $\sim \stackrel{\text{iid}}{\sim} \mathbf{P}_t$ .

# Topic model

## Topic model for single-topic documents



$K$  topics (dists. over words)  $\{\mathbf{P}_t\}_{t=1}^K$ .

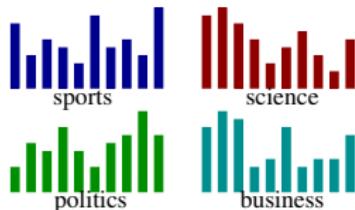
Pick topic  $t$  with prob.  $w_t$  (hidden).

Word tokens in doc.  $\stackrel{\text{iid}}{\sim} \mathbf{P}_t$ .

Given iid sample of documents of length  $L$ ,  
produce estimates of **model parameters**  $\{(\mathbf{P}_t, w_t)\}_{t=1}^K$ .

# Topic model

## Topic model for single-topic documents



$K$  topics (dists. over words)  $\{\mathbf{P}_t\}_{t=1}^K$ .

Pick topic  $t$  with prob.  $w_t$  (hidden).

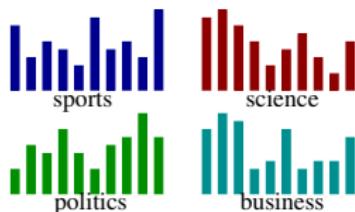
Word tokens in doc.  $\stackrel{\text{iid}}{\sim} \mathbf{P}_t$ .

Given iid sample of documents of length  $L$ ,  
produce estimates of **model parameters**  $\{(\mathbf{P}_t, w_t)\}_{t=1}^K$ .

How long must the documents be?

# Topic model

## Topic model for single-topic documents



$K$  topics (dists. over words)  $\{\mathbf{P}_t\}_{t=1}^K$ .

Pick topic  $t$  with prob.  $w_t$  (hidden).

Word tokens in doc.  $\stackrel{\text{iid}}{\sim} \mathbf{P}_t$ .

Given iid sample of documents of length  $L$ ,  
produce estimates of **model parameters**  $\{(\mathbf{P}_t, w_t)\}_{t=1}^K$ .

How long must the documents be?

(Answering this question leads to efficient algorithms for  
estimating parameters!)

# Identifiability

**Generative process:**

Pick  $t \sim \text{Categorical}(w_1, w_2, \dots, w_K)$ .

Given  $t$ , pick  $L$  words from  $\mathbf{P}_t$ .

# Identifiability

## Generative process:

Pick  $t \sim \text{Categorical}(w_1, w_2, \dots, w_K)$ .

Given  $t$ , pick  $L$  words from  $\mathbf{P}_t$ .

- $L = 1$ : random document (single word)  $\sim \sum_{t=1}^K w_t \mathbf{P}_t$ .

Are parameters  $\{(\mathbf{P}_t, w_t)\}_{t=1}^K$  identifiable from single-word documents?

# Identifiability

## Generative process:

Pick  $t \sim \text{Categorical}(w_1, w_2, \dots, w_K)$ .

Given  $t$ , pick  $L$  words from  $\mathbf{P}_t$ .

- $L = 1$ : random document (single word)  $\sim \sum_{t=1}^K w_t \mathbf{P}_t$ .

Are parameters  $\{(\mathbf{P}_t, w_t)\}_{t=1}^K$  identifiable from single-word documents?

No.

# Identifiability

## Generative process:

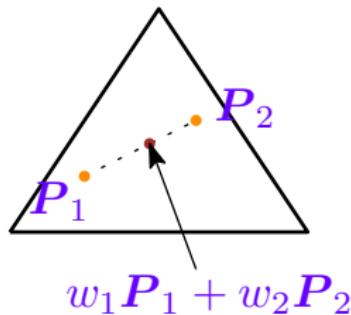
Pick  $t \sim \text{Categorical}(w_1, w_2, \dots, w_K)$ .

Given  $t$ , pick  $L$  words from  $\mathbf{P}_t$ .

- $L = 1$ : random document (single word)  $\sim \sum_{t=1}^K w_t \mathbf{P}_t$ .

Are parameters  $\{(\mathbf{P}_t, w_t)\}_{t=1}^K$  identifiable from single-word documents?

No.



# Identifiability

**Generative process:**

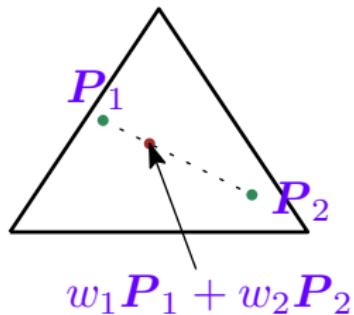
Pick  $t \sim \text{Categorical}(w_1, w_2, \dots, w_K)$ .

Given  $t$ , pick  $L$  words from  $\mathbf{P}_t$ .

- $L = 1$ : random document (single word)  $\sim \sum_{t=1}^K w_t \mathbf{P}_t$ .

Are parameters  $\{(\mathbf{P}_t, w_t)\}_{t=1}^K$  identifiable from single-word documents?

No.



## Identifiability: $L = 2$

**Generative process:**

Pick  $t \sim \text{Categorical}(w_1, w_2, \dots, w_K)$ .

Given  $t$ , pick  $L$  words from  $\mathbf{P}_t$ .

- ▶  $L = 2$ :

# Identifiability: $L = 2$

## Generative process:

Pick  $t \sim \text{Categorical}(w_1, w_2, \dots, w_K)$ .

Given  $t$ , pick  $L$  words from  $\mathbf{P}_t$ .

- $L = 2$ :

Regard  $\mathbf{P}_t$  as probability vector ( $i$ th entry of  $\mathbf{P}_t$  is  $\Pr[\text{word } i]$ ).

Joint distribution of word pairs (for topic  $t$ ) is given by matrix:

$$\mathbf{P}_t \mathbf{P}_t^\top = \begin{matrix} & j \\ i & \begin{array}{|c|c|c|c|} \hline & \cdots & \cdots & \cdots \\ \hline \cdots & \vdots & \vdots & \vdots \\ \hline \cdots & \vdots & \vdots & \vdots \\ \hline \end{array} \end{matrix} \quad \Pr[\text{words } i, j]$$

Random document  $\sim \sum_{t=1}^K w_t \mathbf{P}_t \mathbf{P}_t^\top$ .

# Identifiability: $L = 2$

## Generative process:

Pick  $t \sim \text{Categorical}(w_1, w_2, \dots, w_K)$ .

Given  $t$ , pick  $L$  words from  $\mathbf{P}_t$ .

- ▶  $L = 2$ :

Regard  $\mathbf{P}_t$  as probability vector ( $i$ th entry of  $\mathbf{P}_t$  is  $\Pr[\text{word } i]$ ).

Joint distribution of word pairs (for topic  $t$ ) is given by matrix:

$$\mathbf{P}_t \mathbf{P}_t^\top = \begin{matrix} & j \\ i & \begin{array}{c|c} & \cdots \\ \cdots & \cdots \\ \cdots & \cdots \\ \cdots & \cdots \\ \hline & j \end{array} \end{matrix} \quad \Pr[\text{words } i, j]$$

Random document  $\sim \sum_{t=1}^K w_t \mathbf{P}_t \mathbf{P}_t^\top$ .

Are parameters  $\{(\mathbf{P}_t, w_t)\}_{t=1}^K$  identifiable from word pairs?

## Simple observation

Suppose distribution of word pairs (as a matrix) can be written as

$$\mathbf{M} = \mathbf{A}\mathbf{A}^T.$$

## Simple observation

Suppose distribution of word pairs (as a matrix) can be written as

$$\mathbf{M} = \mathbf{A}\mathbf{A}^T.$$

Then it can also be written as

$$\mathbf{M} = (\mathbf{A}\mathbf{R})(\mathbf{A}\mathbf{R})^T$$

for any orthogonal matrix  $\mathbf{R}$  (because  $\mathbf{R}^T\mathbf{R} = \mathbf{I}$ ).

## Identifiability: $L = 2$ counterexample

Parameters  $\{(\mathbf{P}_1, w_1), (\mathbf{P}_2, w_2)\}$  and  $\{(\tilde{\mathbf{P}}_1, \tilde{w}_1), (\tilde{\mathbf{P}}_2, \tilde{w}_2)\}$

$$(\mathbf{P}_1, w_1) = \left( \begin{bmatrix} 0.40 \\ 0.60 \end{bmatrix}, 0.5 \right), \quad (\mathbf{P}_2, w_2) = \left( \begin{bmatrix} 0.60 \\ 0.40 \end{bmatrix}, 0.5 \right);$$
$$(\tilde{\mathbf{P}}_1, \tilde{w}_1) = \left( \begin{bmatrix} 0.55 \\ 0.45 \end{bmatrix}, 0.8 \right), \quad (\tilde{\mathbf{P}}_2, \tilde{w}_2) = \left( \begin{bmatrix} 0.30 \\ 0.70 \end{bmatrix}, 0.2 \right)$$

## Identifiability: $L = 2$ counterexample

Parameters  $\{(\mathbf{P}_1, w_1), (\mathbf{P}_2, w_2)\}$  and  $\{(\tilde{\mathbf{P}}_1, \tilde{w}_1), (\tilde{\mathbf{P}}_2, \tilde{w}_2)\}$

$$(\mathbf{P}_1, w_1) = \left( \begin{bmatrix} 0.40 \\ 0.60 \end{bmatrix}, 0.5 \right), \quad (\mathbf{P}_2, w_2) = \left( \begin{bmatrix} 0.60 \\ 0.40 \end{bmatrix}, 0.5 \right);$$
$$(\tilde{\mathbf{P}}_1, \tilde{w}_1) = \left( \begin{bmatrix} 0.55 \\ 0.45 \end{bmatrix}, 0.8 \right), \quad (\tilde{\mathbf{P}}_2, \tilde{w}_2) = \left( \begin{bmatrix} 0.30 \\ 0.70 \end{bmatrix}, 0.2 \right)$$

satisfy

$$w_1 \mathbf{P}_1 \mathbf{P}_1^\top + w_2 \mathbf{P}_2 \mathbf{P}_2^\top = \tilde{w}_1 \tilde{\mathbf{P}}_1 \tilde{\mathbf{P}}_1^\top + \tilde{w}_2 \tilde{\mathbf{P}}_2 \tilde{\mathbf{P}}_2^\top = \begin{bmatrix} 0.26 & 0.24 \\ 0.24 & 0.26 \end{bmatrix}.$$

## Identifiability: $L = 2$ counterexample

Parameters  $\{(\mathbf{P}_1, w_1), (\mathbf{P}_2, w_2)\}$  and  $\{(\tilde{\mathbf{P}}_1, \tilde{w}_1), (\tilde{\mathbf{P}}_2, \tilde{w}_2)\}$

$$(\mathbf{P}_1, w_1) = \left( \begin{bmatrix} 0.40 \\ 0.60 \end{bmatrix}, 0.5 \right), \quad (\mathbf{P}_2, w_2) = \left( \begin{bmatrix} 0.60 \\ 0.40 \end{bmatrix}, 0.5 \right);$$
$$(\tilde{\mathbf{P}}_1, \tilde{w}_1) = \left( \begin{bmatrix} 0.55 \\ 0.45 \end{bmatrix}, 0.8 \right), \quad (\tilde{\mathbf{P}}_2, \tilde{w}_2) = \left( \begin{bmatrix} 0.30 \\ 0.70 \end{bmatrix}, 0.2 \right)$$

satisfy

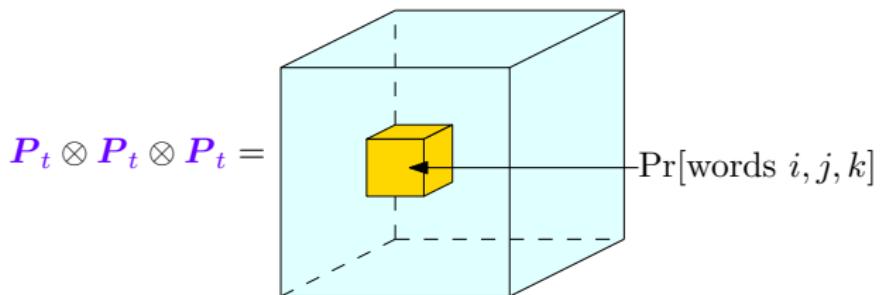
$$w_1 \mathbf{P}_1 \mathbf{P}_1^\top + w_2 \mathbf{P}_2 \mathbf{P}_2^\top = \tilde{w}_1 \tilde{\mathbf{P}}_1 \tilde{\mathbf{P}}_1^\top + \tilde{w}_2 \tilde{\mathbf{P}}_2 \tilde{\mathbf{P}}_2^\top = \begin{bmatrix} 0.26 & 0.24 \\ 0.24 & 0.26 \end{bmatrix}.$$

Cannot identify parameters from length-two documents.

## Identifiability: $L = 3$

**Documents of length  $L = 3$**

Joint distribution of word triple (for topic  $t$ ) is given by *tensor*:



$$\text{Random document} \sim \sum_{t=1}^K w_t \mathbf{P}_t \otimes \mathbf{P}_t \otimes \mathbf{P}_t.$$

## Identifiability from documents of length three

**Claim:** If  $\{\mathbf{P}_t\}_{t=1}^K$  are linearly independent & all  $w_t > 0$ , then parameters  $\{(\mathbf{P}_t, w_t)\}_{t=1}^K$  are identifiable from word triples.

## Identifiability from documents of length three

**Claim:** If  $\{\mathbf{P}_t\}_{t=1}^K$  are linearly independent & all  $w_t > 0$ , then parameters  $\{(\mathbf{P}_t, w_t)\}_{t=1}^K$  are identifiable from word triples.

- ▶ Claim implied by uniqueness of certain *tensor decompositions*.

## Identifiability from documents of length three

**Claim:** If  $\{\mathbf{P}_t\}_{t=1}^K$  are linearly independent & all  $w_t > 0$ , then parameters  $\{(\mathbf{P}_t, w_t)\}_{t=1}^K$  are identifiable from word triples.

- ▶ Claim implied by uniqueness of certain *tensor decompositions*.
- ▶ Proof is *constructive*: i.e., comes with an algorithm!

## Identifiability from documents of length three

**Claim:** If  $\{\mathbf{P}_t\}_{t=1}^K$  are linearly independent & all  $w_t > 0$ , then parameters  $\{(\mathbf{P}_t, w_t)\}_{t=1}^K$  are identifiable from word triples.

- ▶ Claim implied by uniqueness of certain *tensor decompositions*.
- ▶ Proof is *constructive*: i.e., comes with an algorithm!

**Next:** Brief overview of tensors.

## Tensors of order two

**Matrices (tensors of order two):**  $\mathbf{M} \in \mathbb{R}^{d \times d}$ .

- ▶ Regard as *bi-linear function*  $\mathbf{M}: \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$ :

$$\mathbf{M}(ax + x', \mathbf{y}) = a\mathbf{M}(\mathbf{x}, \mathbf{y}) + \mathbf{M}(x', \mathbf{y});$$

$$\mathbf{M}(\mathbf{x}, ay + \mathbf{y}') = a\mathbf{M}(\mathbf{x}, \mathbf{y}) + \mathbf{M}(\mathbf{x}, \mathbf{y}').$$

## Tensors of order two

**Matrices (tensors of order two):**  $\mathbf{M} \in \mathbb{R}^{d \times d}$ .

- ▶ Regard as *bi-linear function*  $\mathbf{M}: \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$ :

$$\mathbf{M}(ax + x', \mathbf{y}) = a\mathbf{M}(\mathbf{x}, \mathbf{y}) + \mathbf{M}(x', \mathbf{y});$$

$$\mathbf{M}(\mathbf{x}, ay + \mathbf{y}') = a\mathbf{M}(\mathbf{x}, \mathbf{y}) + \mathbf{M}(\mathbf{x}, \mathbf{y}').$$

- ▶ Can describe  $\mathbf{M}$  by  $d^2$  values  $\mathbf{M}(e_i, e_j) =: M_{i,j}$ .  
( $e_i$  is  $i$ th coordinate basis vector.)

## Tensors of order two

**Matrices (tensors of order two):**  $\mathbf{M} \in \mathbb{R}^{d \times d}$ .

- ▶ Regard as *bi-linear function*  $\mathbf{M}: \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$ :

$$\begin{aligned}\mathbf{M}(ax + \mathbf{x}', \mathbf{y}) &= a\mathbf{M}(\mathbf{x}, \mathbf{y}) + \mathbf{M}(\mathbf{x}', \mathbf{y}); \\ \mathbf{M}(\mathbf{x}, ay + \mathbf{y}') &= a\mathbf{M}(\mathbf{x}, \mathbf{y}) + \mathbf{M}(\mathbf{x}, \mathbf{y}').\end{aligned}$$

- ▶ Can describe  $\mathbf{M}$  by  $d^2$  values  $\mathbf{M}(e_i, e_j) =: M_{i,j}$ .  
( $e_i$  is  $i$ th coordinate basis vector.)
- ▶ Formula using matrix representation:

$$\mathbf{M}(\mathbf{x}, \mathbf{y}) = \mathbf{x}^\top \mathbf{M} \mathbf{y} = \sum_{i,j} M_{i,j} \cdot x_i y_j.$$

## Tensors of order two

**Matrices (tensors of order two):**  $\mathbf{M} \in \mathbb{R}^{d \times d}$ .

- ▶ Regard as *bi-linear function*  $\mathbf{M}: \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$ :

$$\begin{aligned}\mathbf{M}(ax + \mathbf{x}', \mathbf{y}) &= a\mathbf{M}(\mathbf{x}, \mathbf{y}) + \mathbf{M}(\mathbf{x}', \mathbf{y}); \\ \mathbf{M}(\mathbf{x}, ay + \mathbf{y}') &= a\mathbf{M}(\mathbf{x}, \mathbf{y}) + \mathbf{M}(\mathbf{x}, \mathbf{y}').\end{aligned}$$

- ▶ Can describe  $\mathbf{M}$  by  $d^2$  values  $\mathbf{M}(e_i, e_j) =: M_{i,j}$ .  
( $e_i$  is  $i$ th coordinate basis vector.)
- ▶ Formula using matrix representation:

$$\mathbf{M}(\mathbf{x}, \mathbf{y}) = \mathbf{x}^\top \mathbf{M} \mathbf{y} = \sum_{i,j} M_{i,j} \cdot x_i y_j.$$

Tensors are *multi-linear* generalization.

## Tensors of order $p$

**$p$ -linear functions:**  $\textcolor{violet}{T}: \mathbb{R}^d \times \mathbb{R}^d \times \cdots \times \mathbb{R}^d \rightarrow \mathbb{R}$ .

## Tensors of order $p$

**$p$ -linear functions:**  $\textcolor{violet}{T}: \mathbb{R}^d \times \mathbb{R}^d \times \cdots \times \mathbb{R}^d \rightarrow \mathbb{R}$ .

- ▶ Can describe  $\textcolor{violet}{T}$  by  $d^p$  values  $\textcolor{violet}{T}(e_{i_1}, e_{i_2}, \dots, e_{i_p}) =: \textcolor{violet}{T}_{i_1, i_2, \dots, i_p}$ .

## Tensors of order $p$

**$p$ -linear functions:**  $\mathbf{T}: \mathbb{R}^d \times \mathbb{R}^d \times \cdots \times \mathbb{R}^d \rightarrow \mathbb{R}$ .

- ▶ Can describe  $\mathbf{T}$  by  $d^p$  values  $\mathbf{T}(\mathbf{e}_{i_1}, \mathbf{e}_{i_2}, \dots, \mathbf{e}_{i_p}) =: \mathbf{T}_{i_1, i_2, \dots, i_p}$ .
- ▶ Identify  $\mathbf{T}$  with multi-index array  $\mathbf{T} \in \mathbb{R}^{d \times d \times \cdots \times d}$ .

## Tensors of order $p$

**$p$ -linear functions:**  $\mathbf{T}: \mathbb{R}^d \times \mathbb{R}^d \times \cdots \times \mathbb{R}^d \rightarrow \mathbb{R}$ .

- ▶ Can describe  $\mathbf{T}$  by  $d^p$  values  $\mathbf{T}(e_{i_1}, e_{i_2}, \dots, e_{i_p}) =: \mathbf{T}_{i_1, i_2, \dots, i_p}$ .
- ▶ Identify  $\mathbf{T}$  with multi-index array  $\mathbf{T} \in \mathbb{R}^{d \times d \times \cdots \times d}$ .

Formula for function value:

$$\mathbf{T}(\mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \dots, \mathbf{x}^{(p)}) = \sum_{i_1, i_2, \dots, i_p} T_{i_1, i_2, \dots, i_p} \cdot x_{i_1}^{(1)} x_{i_2}^{(2)} \cdots x_{i_p}^{(p)}.$$

## Tensors of order $p$

**$p$ -linear functions:**  $\mathbf{T}: \mathbb{R}^d \times \mathbb{R}^d \times \cdots \times \mathbb{R}^d \rightarrow \mathbb{R}$ .

- ▶ Can describe  $\mathbf{T}$  by  $d^p$  values  $\mathbf{T}(\mathbf{e}_{i_1}, \mathbf{e}_{i_2}, \dots, \mathbf{e}_{i_p}) =: \mathbf{T}_{i_1, i_2, \dots, i_p}$ .
- ▶ Identify  $\mathbf{T}$  with multi-index array  $\mathbf{T} \in \mathbb{R}^{d \times d \times \cdots \times d}$ .

Formula for function value:

$$\mathbf{T}(\mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \dots, \mathbf{x}^{(p)}) = \sum_{i_1, i_2, \dots, i_p} T_{i_1, i_2, \dots, i_p} \cdot x_{i_1}^{(1)} x_{i_2}^{(2)} \cdots x_{i_p}^{(p)}.$$

- ▶ *Rank-1 tensor:*  $\mathbf{T} = \mathbf{v}^{(1)} \otimes \mathbf{v}^{(2)} \otimes \cdots \otimes \mathbf{v}^{(p)}$ ,
- $$\mathbf{T}(\mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \dots, \mathbf{x}^{(p)}) = \langle \mathbf{v}^{(1)}, \mathbf{x}^{(1)} \rangle \langle \mathbf{v}^{(2)}, \mathbf{x}^{(2)} \rangle \cdots \langle \mathbf{v}^{(p)}, \mathbf{x}^{(p)} \rangle.$$

## Tensors of order $p$

**$p$ -linear functions:**  $\mathbf{T}: \mathbb{R}^d \times \mathbb{R}^d \times \cdots \times \mathbb{R}^d \rightarrow \mathbb{R}$ .

- ▶ Can describe  $\mathbf{T}$  by  $d^p$  values  $\mathbf{T}(\mathbf{e}_{i_1}, \mathbf{e}_{i_2}, \dots, \mathbf{e}_{i_p}) =: \mathbf{T}_{i_1, i_2, \dots, i_p}$ .
- ▶ Identify  $\mathbf{T}$  with multi-index array  $\mathbf{T} \in \mathbb{R}^{d \times d \times \cdots \times d}$ .

Formula for function value:

$$\mathbf{T}(\mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \dots, \mathbf{x}^{(p)}) = \sum_{i_1, i_2, \dots, i_p} T_{i_1, i_2, \dots, i_p} \cdot x_{i_1}^{(1)} x_{i_2}^{(2)} \cdots x_{i_p}^{(p)}.$$

- ▶ *Rank-1 tensor:*  $\mathbf{T} = \mathbf{v}^{(1)} \otimes \mathbf{v}^{(2)} \otimes \cdots \otimes \mathbf{v}^{(p)}$ ,
- $$\mathbf{T}(\mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \dots, \mathbf{x}^{(p)}) = \langle \mathbf{v}^{(1)}, \mathbf{x}^{(1)} \rangle \langle \mathbf{v}^{(2)}, \mathbf{x}^{(2)} \rangle \cdots \langle \mathbf{v}^{(p)}, \mathbf{x}^{(p)} \rangle.$$

*Symmetric rank-1 tensor:*  $\mathbf{T} = \mathbf{v}^{\otimes p} = \mathbf{v} \otimes \mathbf{v} \otimes \cdots \otimes \mathbf{v}$ ,

$$\mathbf{T}(\mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \dots, \mathbf{x}^{(p)}) = \langle \mathbf{v}, \mathbf{x}^{(1)} \rangle \langle \mathbf{v}, \mathbf{x}^{(2)} \rangle \cdots \langle \mathbf{v}, \mathbf{x}^{(p)} \rangle.$$

# Usual caveat

(Hillar & Lim, 2013)

## Most Tensor Problems Are NP-Hard

CHRISTOPHER J. HILLAR, Mathematical Sciences Research Institute  
LEK-HENG LIM, University of Chicago

We prove that multilinear (tensor) analogues of many efficiently computable problems in numerical linear algebra are NP-hard. Our list includes: determining the feasibility of a system of bilinear equations, deciding whether a 3-tensor possesses a given eigenvalue, singular value, or spectral norm; approximating an eigenvalue, eigenvector, singular vector, or the spectral norm; and determining the rank or best rank-1 approximation of a 3-tensor. Furthermore, we show that restricting these problems to symmetric tensors does not alleviate their NP-hardness. We also explain how deciding nonnegative definiteness of a symmetric 4-tensor is NP-hard and how computing the combinatorial hyperdeterminant is NP-, #P-, and VNP-hard.

## Example: rank

- ▶ Rank of  $\mathbf{T}$ : smallest  $r$  s.t.  $\mathbf{T}$  is sum of  $r$  rank-1 tensors.

## Example: rank

- ▶ Rank of  $\mathbf{T}$ : smallest  $r$  s.t.  $\mathbf{T}$  is sum of  $r$  rank-1 tensors.  
(Computing this is NP-hard.)

## Example: rank

- ▶ Rank of  $\mathbf{T}$ : smallest  $r$  s.t.  $\mathbf{T}$  is sum of  $r$  rank-1 tensors.  
(Computing this is NP-hard.)
- ▶ “Border rank” of  $\mathbf{T}$ : smallest  $r$  s.t. there exists sequence  $(\mathbf{T}_k)_{k \in \mathbb{N}}$  of rank- $r$  tensors with  $\lim_{k \rightarrow \infty} \mathbf{T}_k = \mathbf{T}$ .

## Example: rank

- ▶ Rank of  $\mathbf{T}$ : smallest  $r$  s.t.  $\mathbf{T}$  is sum of  $r$  rank-1 tensors.  
(Computing this is NP-hard.)
- ▶ “Border rank” of  $\mathbf{T}$ : smallest  $r$  s.t. there exists sequence  $(\mathbf{T}_k)_{k \in \mathbb{N}}$  of rank- $r$  tensors with  $\lim_{k \rightarrow \infty} \mathbf{T}_k = \mathbf{T}$ .
- ▶ Rank is not same as border rank!

Define

$$\mathbf{T} := \mathbf{x} \otimes \mathbf{x} \otimes \mathbf{y} + \mathbf{x} \otimes \mathbf{y} \otimes \mathbf{x} + \mathbf{y} \otimes \mathbf{x} \otimes \mathbf{x},$$

which has rank 3.

## Example: rank

- ▶ Rank of  $\mathbf{T}$ : smallest  $r$  s.t.  $\mathbf{T}$  is sum of  $r$  rank-1 tensors.  
(Computing this is NP-hard.)
- ▶ “Border rank” of  $\mathbf{T}$ : smallest  $r$  s.t. there exists sequence  $(\mathbf{T}_k)_{k \in \mathbb{N}}$  of rank- $r$  tensors with  $\lim_{k \rightarrow \infty} \mathbf{T}_k = \mathbf{T}$ .
- ▶ Rank is not same as border rank!

Define

$$\mathbf{T} := \mathbf{x} \otimes \mathbf{x} \otimes \mathbf{y} + \mathbf{x} \otimes \mathbf{y} \otimes \mathbf{x} + \mathbf{y} \otimes \mathbf{x} \otimes \mathbf{x},$$

which has rank 3.

Define

$$\mathbf{T}_{1/\epsilon} := \frac{1}{\epsilon}(\mathbf{x} + \epsilon \mathbf{y}) \otimes (\mathbf{x} + \epsilon \mathbf{y}) \otimes (\mathbf{x} + \epsilon \mathbf{y}) - \frac{1}{\epsilon} \mathbf{x} \otimes \mathbf{x} \otimes \mathbf{x},$$

which have rank 2.

## Example: rank

- ▶ Rank of  $\mathbf{T}$ : smallest  $r$  s.t.  $\mathbf{T}$  is sum of  $r$  rank-1 tensors.  
(Computing this is NP-hard.)
- ▶ “Border rank” of  $\mathbf{T}$ : smallest  $r$  s.t. there exists sequence  $(\mathbf{T}_k)_{k \in \mathbb{N}}$  of rank- $r$  tensors with  $\lim_{k \rightarrow \infty} \mathbf{T}_k = \mathbf{T}$ .
- ▶ Rank is not same as border rank!

Define

$$\mathbf{T} := \mathbf{x} \otimes \mathbf{x} \otimes \mathbf{y} + \mathbf{x} \otimes \mathbf{y} \otimes \mathbf{x} + \mathbf{y} \otimes \mathbf{x} \otimes \mathbf{x},$$

which has rank 3.

Define

$$\mathbf{T}_{1/\epsilon} := \frac{1}{\epsilon} (\mathbf{x} + \epsilon \mathbf{y}) \otimes (\mathbf{x} + \epsilon \mathbf{y}) \otimes (\mathbf{x} + \epsilon \mathbf{y}) - \frac{1}{\epsilon} \mathbf{x} \otimes \mathbf{x} \otimes \mathbf{x},$$

which have rank 2.

For  $\epsilon = 1/k$ , have  $\lim_{k \rightarrow \infty} \mathbf{T}_k = \mathbf{T}$ .

## Aside: eigenvalue decomposition

**Recall:** every symmetric matrix  $\mathbf{M} \in \mathbb{R}^{d \times d}$  of rank  $K$  has an *eigen-decomposition* (which can be efficiently computed):

$$\mathbf{M} = \sum_{t=1}^K \lambda_t \mathbf{v}_t \mathbf{v}_t^\top,$$

## Aside: eigenvalue decomposition

**Recall:** every symmetric matrix  $\mathbf{M} \in \mathbb{R}^{d \times d}$  of rank  $K$  has an *eigen-decomposition* (which can be efficiently computed):

$$\mathbf{M} = \sum_{t=1}^K \lambda_t \mathbf{v}_t \mathbf{v}_t^\top,$$

- ▶  $\{\lambda_t\}_{t=1}^K$  are *eigenvalues*,
- ▶  $\{\mathbf{v}_t\}_{t=1}^K$  are the corresponding *eigenvectors*, which are orthonormal (i.e., orthogonal & unit length).
- ▶ Decomposition is unique iff  $\{\lambda_t\}_{t=1}^K$  are distinct.  
(Up to sign of  $\mathbf{v}_t$ s.)

## Aside: eigenvalue decomposition

**Recall:** every symmetric matrix  $\mathbf{M} \in \mathbb{R}^{d \times d}$  of rank  $K$  has an *eigen-decomposition* (which can be efficiently computed):

$$\mathbf{M} = \sum_{t=1}^K \lambda_t \mathbf{v}_t \mathbf{v}_t^\top,$$

- ▶  $\{\lambda_t\}_{t=1}^K$  are *eigenvalues*,
- ▶  $\{\mathbf{v}_t\}_{t=1}^K$  are the corresponding *eigenvectors*, which are orthonormal (i.e., orthogonal & unit length).
- ▶ Decomposition is unique iff  $\{\lambda_t\}_{t=1}^K$  are distinct.  
(Up to sign of  $\mathbf{v}_t$ s.)

For (symmetric) tensors of order  $p \geq 3$ :

an analogous decomposition is **not** guaranteed to exist.

## Reduction to orthonormal case

Suppose we have (estimates of) moments of the form

$$\mathbf{M} = \sum_{t=1}^K \mathbf{v}_t \otimes \mathbf{v}_t, \quad (\text{e.g., word pairs})$$

and  $\mathbf{T} = \sum_{t=1}^K \lambda_t \cdot \mathbf{v}_t \otimes \mathbf{v}_t \otimes \mathbf{v}_t.$  (e.g., word triples)

Here, we assume  $\{\mathbf{v}_t\}_{t=1}^K$  are linearly independent, and  $\{\lambda_t\}_{t=1}^K$  are positive.

## Reduction to orthonormal case

Suppose we have (estimates of) moments of the form

$$\mathbf{M} = \sum_{t=1}^K \mathbf{v}_t \otimes \mathbf{v}_t, \quad (\text{e.g., word pairs})$$

and  $\mathbf{T} = \sum_{t=1}^K \lambda_t \cdot \mathbf{v}_t \otimes \mathbf{v}_t \otimes \mathbf{v}_t.$  (e.g., word triples)

Here, we assume  $\{\mathbf{v}_t\}_{t=1}^K$  are linearly independent, and  $\{\lambda_t\}_{t=1}^K$  are positive.

- ▶  $\mathbf{M}$  is positive semidefinite of rank  $K.$

## Reduction to orthonormal case

Suppose we have (estimates of) moments of the form

$$\mathbf{M} = \sum_{t=1}^K \mathbf{v}_t \otimes \mathbf{v}_t, \quad (\text{e.g., word pairs})$$

and  $\mathbf{T} = \sum_{t=1}^K \lambda_t \cdot \mathbf{v}_t \otimes \mathbf{v}_t \otimes \mathbf{v}_t.$  (e.g., word triples)

Here, we assume  $\{\mathbf{v}_t\}_{t=1}^K$  are linearly independent, and  $\{\lambda_t\}_{t=1}^K$  are positive.

- ▶  $\mathbf{M}$  is positive semidefinite of rank  $K.$
- ▶  $\mathbf{M}$  determines inner product system on  $\text{span}\{\mathbf{v}_t\}_{t=1}^K$  s.t.  
 $\{\mathbf{v}_t\}_{t=1}^K$  are **orthonormal**:

$$\langle \mathbf{x}, \mathbf{y} \rangle_{\mathbf{M}} := \mathbf{x}^\top \mathbf{M}^\dagger \mathbf{y}.$$

## Reduction to orthonormal case

Suppose we have (estimates of) moments of the form

$$\mathbf{M} = \sum_{t=1}^K \mathbf{v}_t \otimes \mathbf{v}_t, \quad (\text{e.g., word pairs})$$

and  $\mathbf{T} = \sum_{t=1}^K \lambda_t \cdot \mathbf{v}_t \otimes \mathbf{v}_t \otimes \mathbf{v}_t.$  (e.g., word triples)

Here, we assume  $\{\mathbf{v}_t\}_{t=1}^K$  are linearly independent, and  $\{\lambda_t\}_{t=1}^K$  are positive.

- ▶  $\mathbf{M}$  is positive semidefinite of rank  $K.$
- ▶  $\mathbf{M}$  determines inner product system on  $\text{span}\{\mathbf{v}_t\}_{t=1}^K$  s.t.  
 $\{\mathbf{v}_t\}_{t=1}^K$  are **orthonormal**:

$$\langle \mathbf{x}, \mathbf{y} \rangle_{\mathbf{M}} := \mathbf{x}^\top \mathbf{M}^\dagger \mathbf{y}.$$

- ▶ ∴ Can assume  $d = K$  and  $\{\mathbf{v}_t\}_{t=1}^d$  are orthonormal.  
(Similar to PCA; called “whitening” in signal processing context.)

## Orthogonally decomposable tensors ( $d = K$ )

**Goal:** Given tensor  $\textcolor{teal}{T} = \sum_{t=1}^d \lambda_t \cdot \mathbf{v}_t \otimes \mathbf{v}_t \otimes \mathbf{v}_t \in \mathbb{R}^{d \times d \times d}$   
where:

- ▶  $\{\mathbf{v}_t\}_{t=1}^d$  are orthonormal;
- ▶ all  $\lambda_t > 0$ ;

approximately recover  $\{(\mathbf{v}_t, \lambda_t)\}_{t=1}^d$ .

## Exact orthogonally decomposable tensor

(Zhang & Golub, 2001)

**Matching moments:**

$$\{(\hat{\mathbf{v}}_t, \hat{\lambda}_t)\}_{t=1}^d := \arg \min_{\{(\mathbf{x}_t, \sigma_t)\}_{t=1}^d} \left\| \mathbf{T} - \sum_{t=1}^d \sigma_t \cdot \mathbf{x}_t \otimes \mathbf{x}_t \otimes \mathbf{x}_t \right\|_F^2.$$

(Here,  $\|\cdot\|_F$  is “Frobenius norm”, just like for matrices.)

# Exact orthogonally decomposable tensor

(Zhang & Golub, 2001)

**Matching moments:**

$$\{(\hat{\mathbf{v}}_t, \hat{\lambda}_t)\}_{t=1}^d := \arg \min_{\{(\mathbf{x}_t, \sigma_t)\}_{t=1}^d} \left\| \mathbf{T} - \sum_{t=1}^d \sigma_t \cdot \mathbf{x}_t \otimes \mathbf{x}_t \otimes \mathbf{x}_t \right\|_F^2.$$

(Here,  $\|\cdot\|_F$  is “Frobenius norm”, just like for matrices.)

- ▶ Greedy approach:
  - ▶ Find best rank-1 approximation:

$$(\hat{\mathbf{v}}, \hat{\lambda}) := \arg \min_{\|\mathbf{x}\|=1, \sigma \geq 0} \|\mathbf{T} - \sigma \cdot \mathbf{x} \otimes \mathbf{x} \otimes \mathbf{x}\|_F^2.$$

- ▶ “Deflate”  $\mathbf{T} := \mathbf{T} - \hat{\lambda} \cdot \hat{\mathbf{v}} \otimes \hat{\mathbf{v}} \otimes \hat{\mathbf{v}}$  and repeat.

# Exact orthogonally decomposable tensor

(Zhang & Golub, 2001)

**Matching moments:**

$$\{(\hat{\mathbf{v}}_t, \hat{\lambda}_t)\}_{t=1}^d := \arg \min_{\{(\mathbf{x}_t, \sigma_t)\}_{t=1}^d} \left\| \mathbf{T} - \sum_{t=1}^d \sigma_t \cdot \mathbf{x}_t \otimes \mathbf{x}_t \otimes \mathbf{x}_t \right\|_F^2.$$

(Here,  $\|\cdot\|_F$  is “Frobenius norm”, just like for matrices.)

- ▶ Greedy approach:
  - ▶ Find best rank-1 approximation:

$$\hat{\mathbf{v}} := \arg \max_{\|\mathbf{x}\|=1} \mathbf{T}(\mathbf{x}, \mathbf{x}, \mathbf{x}), \quad \hat{\lambda} := \mathbf{T}(\hat{\mathbf{v}}, \hat{\mathbf{v}}, \hat{\mathbf{v}}).$$

- ▶ “Deflate”  $\mathbf{T} := \mathbf{T} - \hat{\lambda} \cdot \hat{\mathbf{v}} \otimes \hat{\mathbf{v}} \otimes \hat{\mathbf{v}}$  and repeat.

## Rank-1 approximation problem

**Claim:** Local maximizers of the function

$$\mathbf{x} \mapsto \mathbf{T}(\mathbf{x}, \mathbf{x}, \mathbf{x}) = \sum_{i,j,k} \mathbf{T}_{i,j,k} \cdot x_i x_j x_k$$

(over the unit ball) are  $\{\mathbf{v}_t\}_{t=1}^d$ , and

$$\mathbf{T}(\mathbf{v}_t, \mathbf{v}_t, \mathbf{v}_t) = \lambda_t, \quad t \in [d].$$

## Rank-1 approximation problem

**Claim:** Local maximizers of the function

$$\mathbf{x} \mapsto \mathbf{T}(\mathbf{x}, \mathbf{x}, \mathbf{x}) = \sum_{i,j,k} T_{i,j,k} \cdot x_i x_j x_k = \sum_{t=1}^d \lambda_t \cdot \langle \mathbf{v}_t, \mathbf{x} \rangle^3$$

(over the unit ball) are  $\{\mathbf{v}_t\}_{t=1}^d$ , and

$$\mathbf{T}(\mathbf{v}_t, \mathbf{v}_t, \mathbf{v}_t) = \lambda_t, \quad t \in [d].$$

## Rank-1 approximation problem

**Claim:** Local maximizers of the function

$$\mathbf{x} \mapsto \mathbf{T}(\mathbf{x}, \mathbf{x}, \mathbf{x}) = \sum_{i,j,k} \mathbf{T}_{i,j,k} \cdot x_i x_j x_k = \sum_{t=1}^d \lambda_t \cdot \langle \mathbf{v}_t, \mathbf{x} \rangle^3$$

(over the unit ball) are  $\{\mathbf{v}_t\}_{t=1}^d$ , and

$$\mathbf{T}(\mathbf{v}_t, \mathbf{v}_t, \mathbf{v}_t) = \lambda_t, \quad t \in [d].$$

**Corollary:** decomposition of  $\mathbf{T}$  as  $\sum_{t=1}^K \lambda_t \cdot \mathbf{v}_t^{\otimes 3}$  is *unique*!

## Proof

By linearity and orthogonality:

$$T(\mathbf{v}_t, \mathbf{v}_t, \mathbf{v}_t) = \sum_{s=1}^d (\lambda_s \cdot \mathbf{v}_s^{\otimes 3})(\mathbf{v}_t, \mathbf{v}_t, \mathbf{v}_t)$$

## Proof

By linearity and orthogonality:

$$T(\mathbf{v}_t, \mathbf{v}_t, \mathbf{v}_t) = \sum_{s=1}^d (\lambda_s \cdot \mathbf{v}_s^{\otimes 3})(\mathbf{v}_t, \mathbf{v}_t, \mathbf{v}_t) = \sum_{s=1}^d \begin{cases} \lambda_s & \text{if } s = t \\ 0 & \text{if } s \neq t \end{cases}$$

## Proof

By linearity and orthogonality:

$$T(\mathbf{v}_t, \mathbf{v}_t, \mathbf{v}_t) = \sum_{s=1}^d (\lambda_s \cdot \mathbf{v}_s^{\otimes 3})(\mathbf{v}_t, \mathbf{v}_t, \mathbf{v}_t) = \sum_{s=1}^d \begin{cases} \lambda_s & \text{if } s = t \\ 0 & \text{if } s \neq t \end{cases} = \lambda_t.$$

## Proof

By linearity and orthogonality:

$$\mathbf{T}(\mathbf{v}_t, \mathbf{v}_t, \mathbf{v}_t) = \sum_{s=1}^d (\lambda_s \cdot \mathbf{v}_s^{\otimes 3})(\mathbf{v}_t, \mathbf{v}_t, \mathbf{v}_t) = \sum_{s=1}^d \begin{cases} \lambda_s & \text{if } s = t \\ 0 & \text{if } s \neq t \end{cases} = \lambda_t.$$

WLOG assume  $\mathbf{v}_t = e_t$ , so optimization problem is

$$\max_{x \in \mathbb{R}^d} \sum_{t=1}^d \lambda_t x_t^3 \quad \text{s.t.} \quad \sum_{t=1}^d x_t^2 \leq 1.$$

## Proof

By linearity and orthogonality:

$$\mathbf{T}(\mathbf{v}_t, \mathbf{v}_t, \mathbf{v}_t) = \sum_{s=1}^d (\lambda_s \cdot \mathbf{v}_s^{\otimes 3})(\mathbf{v}_t, \mathbf{v}_t, \mathbf{v}_t) = \sum_{s=1}^d \begin{cases} \lambda_s & \text{if } s = t \\ 0 & \text{if } s \neq t \end{cases} = \lambda_t.$$

WLOG assume  $\mathbf{v}_t = \mathbf{e}_t$ , so optimization problem is

$$\max_{\mathbf{x} \in \mathbb{R}^d} \sum_{t=1}^d \lambda_t x_t^3 \quad \text{s.t.} \quad \sum_{t=1}^d x_t^2 \leq 1.$$

If both  $x_1$  and  $x_2$  are non-zero, then

$$\lambda_1 x_1^3 + \lambda_2 x_2^3 < \lambda_1 x_1^2 + \lambda_2 x_2^2$$

## Proof

By linearity and orthogonality:

$$\mathbf{T}(\mathbf{v}_t, \mathbf{v}_t, \mathbf{v}_t) = \sum_{s=1}^d (\lambda_s \cdot \mathbf{v}_s^{\otimes 3})(\mathbf{v}_t, \mathbf{v}_t, \mathbf{v}_t) = \sum_{s=1}^d \begin{cases} \lambda_s & \text{if } s = t \\ 0 & \text{if } s \neq t \end{cases} = \lambda_t.$$

WLOG assume  $\mathbf{v}_t = \mathbf{e}_t$ , so optimization problem is

$$\max_{x \in \mathbb{R}^d} \sum_{t=1}^d \lambda_t x_t^3 \quad \text{s.t.} \quad \sum_{t=1}^d x_t^2 \leq 1.$$

If both  $x_1$  and  $x_2$  are non-zero, then

$$\lambda_1 x_1^3 + \lambda_2 x_2^3 < \lambda_1 x_1^2 + \lambda_2 x_2^2 \leq \max\{\lambda_1, \lambda_2\}.$$

## Proof

By linearity and orthogonality:

$$\mathbf{T}(\mathbf{v}_t, \mathbf{v}_t, \mathbf{v}_t) = \sum_{s=1}^d (\lambda_s \cdot \mathbf{v}_s^{\otimes 3})(\mathbf{v}_t, \mathbf{v}_t, \mathbf{v}_t) = \sum_{s=1}^d \begin{cases} \lambda_s & \text{if } s = t \\ 0 & \text{if } s \neq t \end{cases} = \lambda_t.$$

WLOG assume  $\mathbf{v}_t = \mathbf{e}_t$ , so optimization problem is

$$\max_{\mathbf{x} \in \mathbb{R}^d} \sum_{t=1}^d \lambda_t x_t^3 \quad \text{s.t.} \quad \sum_{t=1}^d x_t^2 \leq 1.$$

If both  $x_1$  and  $x_2$  are non-zero, then

$$\lambda_1 x_1^3 + \lambda_2 x_2^3 < \lambda_1 x_1^2 + \lambda_2 x_2^2 \leq \max\{\lambda_1, \lambda_2\}.$$

So better to put all energy on a single coordinate.

$\therefore$  Local maximizers are  $\mathbf{e}_1, \mathbf{e}_2, \dots, \mathbf{e}_d$ .

## Proof

By linearity and orthogonality:

$$\mathbf{T}(\mathbf{v}_t, \mathbf{v}_t, \mathbf{v}_t) = \sum_{s=1}^d (\lambda_s \cdot \mathbf{v}_s^{\otimes 3})(\mathbf{v}_t, \mathbf{v}_t, \mathbf{v}_t) = \sum_{s=1}^d \begin{cases} \lambda_s & \text{if } s = t \\ 0 & \text{if } s \neq t \end{cases} = \lambda_t.$$

WLOG assume  $\mathbf{v}_t = e_t$ , so optimization problem is

$$\max_{\mathbf{x} \in \mathbb{R}^d} \sum_{t=1}^d \lambda_t x_t^3 \quad \text{s.t.} \quad \sum_{t=1}^d x_t^2 \leq 1.$$

If both  $x_1$  and  $x_2$  are non-zero, then

$$\lambda_1 x_1^3 + \lambda_2 x_2^3 < \lambda_1 x_1^2 + \lambda_2 x_2^2 \leq \max\{\lambda_1, \lambda_2\}.$$

So better to put all energy on a single coordinate.

∴ Local maximizers are  $\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_d$ . □

# Uniqueness of orthogonal decompositions

## What we have seen so far:

- When components  $\{\mathbf{v}_t\}_{t=1}^d$  are linearly independent:
  - Reduce decomposition problem to *orthogonal tensor decomposition*, where components are orthonormal.

# Uniqueness of orthogonal decompositions

## What we have seen so far:

- When components  $\{\mathbf{v}_t\}_{t=1}^d$  are linearly independent:
  - Reduce decomposition problem to *orthogonal tensor decomposition*, where components are orthonormal.
- For orthogonally decomposable tensors  $\mathbf{T}$ , local maximizers of the function

$$\mathbf{x} \mapsto \mathbf{T}(\mathbf{x}, \mathbf{x}, \mathbf{x})$$

(over the unit ball) are  $\{\mathbf{v}_t\}_{t=1}^d$ .

# Uniqueness of orthogonal decompositions

## What we have seen so far:

- When components  $\{\mathbf{v}_t\}_{t=1}^d$  are linearly independent:
  - Reduce decomposition problem to *orthogonal tensor decomposition*, where components are orthonormal.
- For orthogonally decomposable tensors  $\mathbf{T}$ , local maximizers of the function

$$\mathbf{x} \mapsto \mathbf{T}(\mathbf{x}, \mathbf{x}, \mathbf{x})$$

(over the unit ball) are  $\{\mathbf{v}_t\}_{t=1}^d$ .

**Algorithm:** use gradient ascent to find all of the local maximizers, which are exactly  $\mathbf{v}_t$ .

(Can use “deflation” to remove components from  $\mathbf{T}$  that you’ve already found.)

## Application to topic model parameters

Probabilities of word triples as third-order tensor:

$$\mathbf{T} = \sum_{t=1}^K w_t \mathbf{P}_t \otimes \mathbf{P}_t \otimes \mathbf{P}_t = \sum_{t=1}^K \mathbf{v}_t \otimes \mathbf{v}_t \otimes \mathbf{v}_t$$

for  $\mathbf{v}_t = w_t^{1/3} \mathbf{P}_t$ .

## Application to topic model parameters

Probabilities of word triples as third-order tensor:

$$\mathbf{T} = \sum_{t=1}^K w_t \mathbf{P}_t \otimes \mathbf{P}_t \otimes \mathbf{P}_t = \sum_{t=1}^K \mathbf{v}_t \otimes \mathbf{v}_t \otimes \mathbf{v}_t$$

for  $\mathbf{v}_t = w_t^{1/3} \mathbf{P}_t$ .

- ▶ About linear independence condition on  $\{\mathbf{v}_t\}_{t=1}^K$ :

$\{\mathbf{v}_t\}_{t=1}^K$  are linearly independent

$\Leftrightarrow \{\mathbf{P}_t\}_{t=1}^K$  are linearly independent and all  $w_t > 0$ .

## Application to topic model parameters

Probabilities of word triples as third-order tensor:

$$\mathbf{T} = \sum_{t=1}^K w_t \mathbf{P}_t \otimes \mathbf{P}_t \otimes \mathbf{P}_t = \sum_{t=1}^K \mathbf{v}_t \otimes \mathbf{v}_t \otimes \mathbf{v}_t$$

for  $\mathbf{v}_t = w_t^{1/3} \mathbf{P}_t$ .

- ▶ About linear independence condition on  $\{\mathbf{v}_t\}_{t=1}^K$ :

$\{\mathbf{v}_t\}_{t=1}^K$  are linearly independent

$\Leftrightarrow \{\mathbf{P}_t\}_{t=1}^K$  are linearly independent and all  $w_t > 0$ .

- ▶ Can recover  $\{\mathbf{P}_t\}_{t=1}^K$  from  $\{c_t \mathbf{v}_t\}_{t=1}^K$  for any  $c_t \neq 0$ .

## Application to topic model parameters

Probabilities of word triples as third-order tensor:

$$\mathbf{T} = \sum_{t=1}^K w_t \mathbf{P}_t \otimes \mathbf{P}_t \otimes \mathbf{P}_t = \sum_{t=1}^K \mathbf{v}_t \otimes \mathbf{v}_t \otimes \mathbf{v}_t$$

for  $\mathbf{v}_t = w_t^{1/3} \mathbf{P}_t$ .

- ▶ About linear independence condition on  $\{\mathbf{v}_t\}_{t=1}^K$ :

$\{\mathbf{v}_t\}_{t=1}^K$  are linearly independent

$\Leftrightarrow \{\mathbf{P}_t\}_{t=1}^K$  are linearly independent and all  $w_t > 0$ .

- ▶ Can recover  $\{\mathbf{P}_t\}_{t=1}^K$  from  $\{c_t \mathbf{v}_t\}_{t=1}^K$  for any  $c_t \neq 0$ .
- ▶ Can recover  $\{(\mathbf{P}_t, w_t)\}_{t=1}^K$  from  $\{\mathbf{P}_t\}_{t=1}^K$  and  $\mathbf{T}$ .

□

## Recap

- ▶ Parameters of topic model for single-topic documents (satisfying linear independence condition) can be efficiently recovered from distribution of three-word documents.

## Recap

- ▶ Parameters of topic model for single-topic documents (satisfying linear independence condition) can be efficiently recovered from distribution of three-word documents.
- ▶ Two-word documents not sufficient (without further assumptions).

## Recap

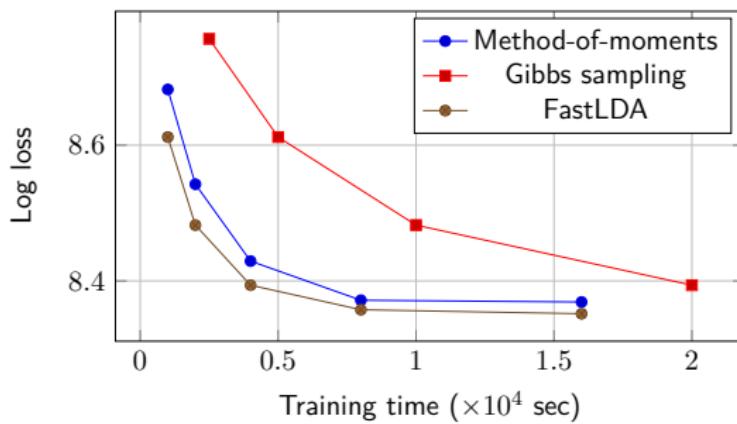
- ▶ Parameters of topic model for single-topic documents (satisfying linear independence condition) can be efficiently recovered from distribution of three-word documents.
- ▶ Two-word documents not sufficient (without further assumptions).
- ▶ Variational characterization of orthogonally decomposable tensors leads to simple and efficient algorithms!

## Illustrative empirical results

- ▶ Corpus: 300,000 New York Times articles.
- ▶ Vocabulary size: 102,660 words.
- ▶ Set number of topics  $K := 50$ .

### Model predictive performance:

≈ 4–8× speed-up over Gibbs sampling for LDA;  
comparable to “FastLDA” (Porteous, Newman, Ihler, Asuncion, Smyth, & Welling, 2008).



## Illustrative empirical results

**Sample topics:** (showing top 10 words for each topic)

Econ.	Baseball	Edu.	Health care	Golf
sales	run	school	drug	player
economic	inning	student	patient	tiger_wood
consumer	hit	teacher	million	won
major	game	program	company	shot
home	season	official	doctor	play
indicator	home	public	companies	round
weekly	right	children	percent	win
order	games	high	cost	tournament
claim	dodger	education	program	tour
scheduled	left	district	health	right

## Illustrative empirical results

**Sample topics:** (showing top 10 words for each topic)

Invest.	Election	auto race	Child's Lit.	Afghan War
percent	al_gore	car	book	taliban
stock	campaign	race	children	attack
market	president	driver	ages	afghanistan
fund	george_bush	team	author	official
investor	bush	won	read	military
companies	clinton	win	newspaper	u_s
analyst	vice	racing	web	united_states
money	presidential	track	writer	terrorist
investment	million	season	written	war
economy	democratic	lap	sales	bin

# Illustrative empirical results

**Sample topics:** (showing top 10 words for each topic)

Web	Antitrust	TV	Movies	Music
com	court	show	film	music
www	case	network	movie	song
site	law	season	director	group
web	lawyer	nbc	play	part
sites	federal	cb	character	new_york
information	government	program	actor	company
online	decision	television	show	million
mail	trial	series	movies	band
internet	microsoft	night	million	show
telegram	right	new_york	part	album

*etc.*

## Computation

**Caveat:** forming and computing with a third-order tensor  $\mathbf{T}$  generally requires **cubic space**.

## Computation

**Caveat:** forming and computing with a third-order tensor  $\mathbf{T}$  generally requires **cubic space**.

- ▶ Fortunately, the tensor we often work with is an *empirical estimate* of a  $\mathbf{T}$ : e.g.,

$$\hat{\mathbf{T}} = \frac{1}{n} \sum_{i=1}^n \text{data}_i,$$

where  $\text{data}_i$  is a tensor involving only the  $i$ -th data point.

## Computation

**Caveat:** forming and computing with a third-order tensor  $\mathbf{T}$  generally requires **cubic space**.

- ▶ Fortunately, the tensor we often work with is an *empirical estimate* of a  $\mathbf{T}$ : e.g.,

$$\hat{\mathbf{T}} = \frac{1}{n} \sum_{i=1}^n \text{data}_i,$$

where  $\text{data}_i$  is a tensor involving only the  $i$ -th data point.

- ▶ Our algorithms will only involve  $\hat{\mathbf{T}}$  through *evaluations* of  $\hat{\mathbf{T}}$  at (several) given arguments, say,  $x, y, z$ .

## Computation

**Caveat:** forming and computing with a third-order tensor  $\mathbf{T}$  generally requires **cubic space**.

- ▶ Fortunately, the tensor we often work with is an *empirical estimate* of a  $\mathbf{T}$ : e.g.,

$$\hat{\mathbf{T}} = \frac{1}{n} \sum_{i=1}^n \text{data}_i,$$

where  $\text{data}_i$  is a tensor involving only the  $i$ -th data point.

- ▶ Our algorithms will only involve  $\hat{\mathbf{T}}$  through *evaluations* of  $\hat{\mathbf{T}}$  at (several) given arguments, say,  $\mathbf{x}, \mathbf{y}, \mathbf{z}$ .

By linearity:

$$\hat{\mathbf{T}}(\mathbf{x}, \mathbf{y}, \mathbf{z}) = \frac{1}{n} \sum_{i=1}^n \text{data}_i(\mathbf{x}, \mathbf{y}, \mathbf{z}).$$

## Computation

**Caveat:** forming and computing with a third-order tensor  $\mathbf{T}$  generally requires **cubic space**.

- ▶ Fortunately, the tensor we often work with is an *empirical estimate* of a  $\mathbf{T}$ : e.g.,

$$\widehat{\mathbf{T}} = \frac{1}{n} \sum_{i=1}^n \text{data}_i,$$

where  $\text{data}_i$  is a tensor involving only the  $i$ -th data point.

- ▶ Our algorithms will only involve  $\widehat{\mathbf{T}}$  through *evaluations* of  $\widehat{\mathbf{T}}$  at (several) given arguments, say,  $\mathbf{x}, \mathbf{y}, \mathbf{z}$ .

By linearity:

$$\widehat{\mathbf{T}}(\mathbf{x}, \mathbf{y}, \mathbf{z}) = \frac{1}{n} \sum_{i=1}^n \text{data}_i(\mathbf{x}, \mathbf{y}, \mathbf{z}).$$

- ▶ **Often:**  $\text{data}_i(\mathbf{x}, \mathbf{y}, \mathbf{z})$  is **easy to compute, even without forming any tensors!**

## Computation

**Caveat:** forming and computing with a third-order tensor  $\mathbf{T}$  generally requires **cubic space**.

- ▶ Fortunately, the tensor we often work with is an *empirical estimate* of a  $\mathbf{T}$ : e.g.,

$$\hat{\mathbf{T}} = \frac{1}{n} \sum_{i=1}^n \text{data}_i,$$

where  $\text{data}_i$  is a tensor involving only the  $i$ -th data point.

- ▶ Our algorithms will only involve  $\hat{\mathbf{T}}$  through *evaluations* of  $\hat{\mathbf{T}}$  at (several) given arguments, say,  $\mathbf{x}, \mathbf{y}, \mathbf{z}$ .

By linearity:

$$\hat{\mathbf{T}}(\mathbf{x}, \mathbf{y}, \mathbf{z}) = \frac{1}{n} \sum_{i=1}^n \text{data}_i(\mathbf{x}, \mathbf{y}, \mathbf{z}).$$

- ▶ **Often:**  $\text{data}_i(\mathbf{x}, \mathbf{y}, \mathbf{z})$  is **easy to compute**, even without forming any tensors! → Linear time/space algorithms.

# Learning algorithms

- ▶ Estimation via **method-of-moments**:
  1. Estimate distribution of three-word documents  $\rightarrow \hat{\mathbf{T}}$  (*empirical moment tensor*).
  2. Approximately decompose  $\hat{\mathbf{T}} \rightarrow$  estimates  $\{(\hat{\mathbf{P}}_t, \hat{w}_t)\}_{t=1}^K$ .

# Learning algorithms

- ▶ Estimation via **method-of-moments**:
  1. Estimate distribution of three-word documents →  $\hat{\mathbf{T}}$  (*empirical moment tensor*).
  2. Approximately decompose  $\hat{\mathbf{T}}$  → estimates  $\{(\hat{\mathbf{P}}_t, \hat{w}_t)\}_{t=1}^K$ .
- ▶ Issues:
  1. Accuracy of *moment estimates*?
  2. Robustness of *(approximate) tensor decomposition*?
  3. *Generality* beyond simple topic models?

# Learning algorithms

- ▶ Estimation via **method-of-moments**:
  1. Estimate distribution of three-word documents →  $\hat{\mathbf{T}}$  (*empirical moment tensor*).
  2. Approximately decompose  $\hat{\mathbf{T}}$  → estimates  $\{(\hat{\mathbf{P}}_t, \hat{w}_t)\}_{t=1}^K$ .
- ▶ Issues:
  1. Accuracy of *moment estimates*?  
Can more reliably estimate lower-order moments;  
distribution-specific sample complexity bounds.
  2. Robustness of *(approximate) tensor decomposition*?  
In some sense, more stable than matrix eigen-decomposition  
(Mu, H., & Goldfarb, 2015)!
  3. *Generality* beyond simple topic models?

# Learning algorithms

- ▶ Estimation via **method-of-moments**:
  1. Estimate distribution of three-word documents →  $\hat{\mathbf{T}}$  (*empirical moment tensor*).
  2. Approximately decompose  $\hat{\mathbf{T}}$  → estimates  $\{(\hat{\mathbf{P}}_t, \hat{w}_t)\}_{t=1}^K$ .
- ▶ Issues:
  1. Accuracy of *moment estimates*?  
Can more reliably estimate lower-order moments;  
distribution-specific sample complexity bounds.
  2. Robustness of *(approximate) tensor decomposition*?  
In some sense, more stable than matrix eigen-decomposition  
(Mu, H., & Goldfarb, 2015)!
  3. *Generality* beyond simple topic models?

**Next:** Moment decompositions for other models.

## 2. Moment decompositions for other models

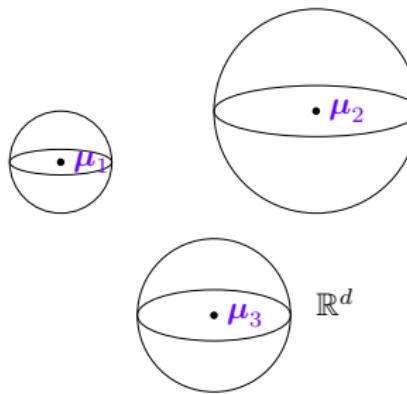
# Moment decompositions

**Some examples of usable moment decompositions.**

1. Two classical mixture models.
2. Models with multi-view structure.
3. Single-index models.

## Mixture model #1: Mixtures of spherical Gaussians

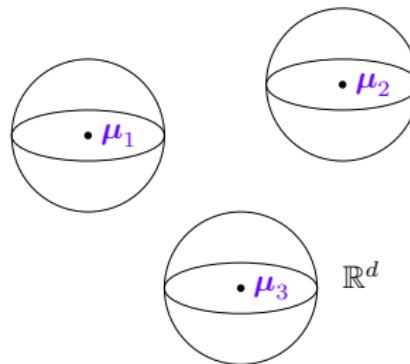
$$H \sim \text{Categorical}(\pi_1, \pi_2, \dots, \pi_K) \quad (\text{hidden});$$
$$\mathbf{X} | H = t \sim \text{Normal}(\boldsymbol{\mu}_t, \sigma_t^2 \mathbf{I}_d), \quad t \in [K].$$



## Mixture model #1: Mixtures of spherical Gaussians

$$H \sim \text{Categorical}(\pi_1, \pi_2, \dots, \pi_K) \quad (\text{hidden});$$
$$\mathbf{X} | H = t \sim \text{Normal}(\boldsymbol{\mu}_t, \sigma^2 \mathbf{I}_d), \quad t \in [K].$$

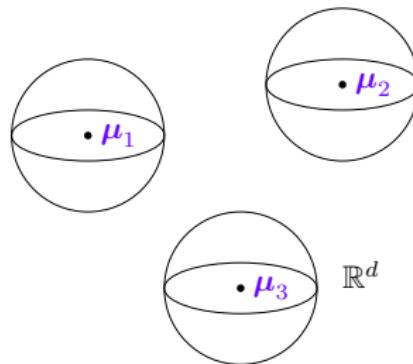
(For simplicity, restrict  $\sigma_1 = \sigma_2 = \dots = \sigma_K = \sigma$ .)



## Mixture model #1: Mixtures of spherical Gaussians

$$H \sim \text{Categorical}(\pi_1, \pi_2, \dots, \pi_K) \quad (\text{hidden});$$
$$\mathbf{X} | H = t \sim \text{Normal}(\boldsymbol{\mu}_t, \sigma^2 \mathbf{I}_d), \quad t \in [K].$$

(For simplicity, restrict  $\sigma_1 = \sigma_2 = \dots = \sigma_K = \sigma$ .)



**Generative process:**

$$\mathbf{X} = \mathbf{Y} + \sigma \mathbf{Z}$$

where  $\Pr(\mathbf{Y} = \boldsymbol{\mu}_t) = \pi_t$ , and  
 $\mathbf{Z} \sim \text{Normal}(\mathbf{0}, \mathbf{I}_d)$  (indep. of  $\mathbf{Y}$ ).

## Using moments for spherical Gaussian mixtures

We'll see two ways to use low-order moments.

## Using moments for spherical Gaussian mixtures

We'll see two ways to use low-order moments.

**First- and second-order moments:**

$$\mathbb{E}(\textcolor{blue}{X}) \in \mathbb{R}^d \quad \text{and} \quad \mathbb{E}(\textcolor{blue}{X} \otimes \textcolor{blue}{X}) \in \mathbb{R}^{d \times d}.$$

## Using moments for spherical Gaussian mixtures

We'll see two ways to use low-order moments.

**First- and second-order moments:**

$$\mathbb{E}(\mathbf{X}) \in \mathbb{R}^d \quad \text{and} \quad \mathbb{E}(\mathbf{X} \otimes \mathbf{X}) \in \mathbb{R}^{d \times d}.$$

**Claim** (Vempala & Wang, 2002):

Span of top  $K$  eigenvectors of  $\mathbb{E}(\mathbf{X} \otimes \mathbf{X})$  contains  $\{\boldsymbol{\mu}_t\}_{t=1}^K$ .

( $K$ -dimensional Principal Component Analysis (PCA) subspace.)

## Proof

**Key fact:**  $k$ -dimensional PCA subspace (based on  $\mathbb{E}(\mathbf{X} \otimes \mathbf{X})$ ) captures as much of overall variance as any other  $k$ -dim. subspace.

## Proof

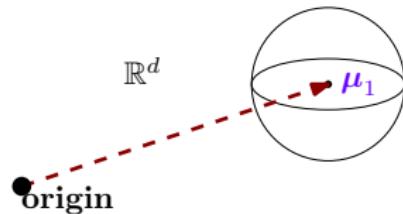
**Key fact:**  $k$ -dimensional PCA subspace (based on  $\mathbb{E}(\mathbf{X} \otimes \mathbf{X})$ ) captures as much of overall variance as any other  $k$ -dim. subspace.

- ▶  $K = 1$  (just a single Gaussian):  
What is the 1-dimensional PCA subspace?

## Proof

**Key fact:**  $k$ -dimensional PCA subspace (based on  $\mathbb{E}(\mathbf{X} \otimes \mathbf{X})$ ) captures as much of overall variance as any other  $k$ -dim. subspace.

- ▶  $K = 1$  (just a single Gaussian):  
What is the 1-dimensional PCA subspace?



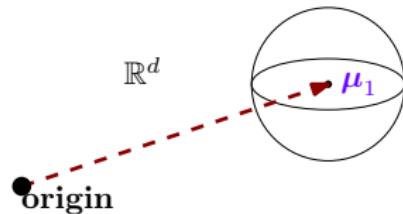
$$\mathbb{E}(\mathbf{X} \otimes \mathbf{X}) = \boldsymbol{\mu}_1 \otimes \boldsymbol{\mu}_1 + \sigma^2 \mathbf{I}_d .$$

## Proof

**Key fact:**  $k$ -dimensional PCA subspace (based on  $\mathbb{E}(\mathbf{X} \otimes \mathbf{X})$ ) captures as much of overall variance as any other  $k$ -dim. subspace.

- ▶  $K = 1$  (just a single Gaussian):

What is the 1-dimensional PCA subspace?



$$\mathbb{E}(\mathbf{X} \otimes \mathbf{X}) = \boldsymbol{\mu}_1 \otimes \boldsymbol{\mu}_1 + \sigma^2 \mathbf{I}_d .$$

Variance in direction  $\mathbf{v}$  (with  $\|\mathbf{v}\| = 1$ ):

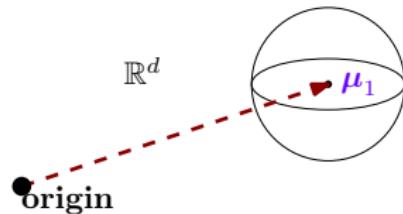
$$\mathbf{v}^\top \mathbb{E}(\mathbf{X} \otimes \mathbf{X}) \mathbf{v}$$

## Proof

**Key fact:**  $k$ -dimensional PCA subspace (based on  $\mathbb{E}(\mathbf{X} \otimes \mathbf{X})$ ) captures as much of overall variance as any other  $k$ -dim. subspace.

- ▶  $K = 1$  (just a single Gaussian):

What is the 1-dimensional PCA subspace?



$$\mathbb{E}(\mathbf{X} \otimes \mathbf{X}) = \boldsymbol{\mu}_1 \otimes \boldsymbol{\mu}_1 + \sigma^2 \mathbf{I}_d .$$

Variance in direction  $\mathbf{v}$  (with  $\|\mathbf{v}\| = 1$ ):

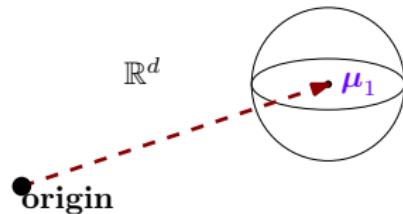
$$\mathbf{v}^\top \mathbb{E}(\mathbf{X} \otimes \mathbf{X}) \mathbf{v} = \mathbf{v}^\top (\boldsymbol{\mu}_1 \otimes \boldsymbol{\mu}_1 + \sigma^2 \mathbf{I}_d) \mathbf{v}$$

## Proof

**Key fact:**  $k$ -dimensional PCA subspace (based on  $\mathbb{E}(\mathbf{X} \otimes \mathbf{X})$ ) captures as much of overall variance as any other  $k$ -dim. subspace.

- ▶  $K = 1$  (just a single Gaussian):

What is the 1-dimensional PCA subspace?



$$\mathbb{E}(\mathbf{X} \otimes \mathbf{X}) = \boldsymbol{\mu}_1 \otimes \boldsymbol{\mu}_1 + \sigma^2 \mathbf{I}_d .$$

Variance in direction  $\mathbf{v}$  (with  $\|\mathbf{v}\| = 1$ ):

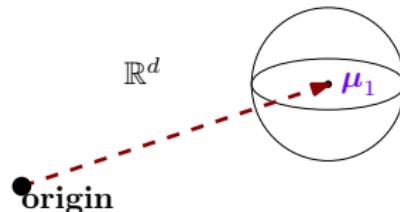
$$\mathbf{v}^\top \mathbb{E}(\mathbf{X} \otimes \mathbf{X}) \mathbf{v} = \mathbf{v}^\top (\boldsymbol{\mu}_1 \otimes \boldsymbol{\mu}_1 + \sigma^2 \mathbf{I}_d) \mathbf{v} = (\mathbf{v}^\top \boldsymbol{\mu}_1)^2 + \sigma^2 .$$

## Proof

**Key fact:**  $k$ -dimensional PCA subspace (based on  $\mathbb{E}(\mathbf{X} \otimes \mathbf{X})$ ) captures as much of overall variance as any other  $k$ -dim. subspace.

- ▶  $K = 1$  (just a single Gaussian):

What is the 1-dimensional PCA subspace?



$$\mathbb{E}(\mathbf{X} \otimes \mathbf{X}) = \boldsymbol{\mu}_1 \otimes \boldsymbol{\mu}_1 + \sigma^2 \mathbf{I}_d .$$

Variance in direction  $\mathbf{v}$  (with  $\|\mathbf{v}\| = 1$ ):

$$\mathbf{v}^\top \mathbb{E}(\mathbf{X} \otimes \mathbf{X}) \mathbf{v} = \mathbf{v}^\top (\boldsymbol{\mu}_1 \otimes \boldsymbol{\mu}_1 + \sigma^2 \mathbf{I}_d) \mathbf{v} = (\mathbf{v}^\top \boldsymbol{\mu}_1)^2 + \sigma^2 .$$

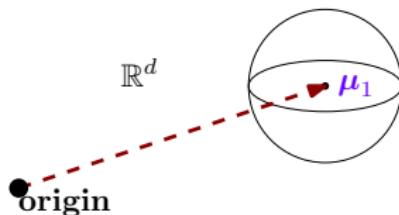
**Best direction** (1-dim. PCA subspace):  $\mathbf{v} = \pm \boldsymbol{\mu}_1 / \|\boldsymbol{\mu}_1\|$ .

## Proof (continued)

**Key fact:**  $k$ -dimensional PCA subspace (based on  $\mathbb{E}(\mathbf{X} \otimes \mathbf{X})$ ) captures as much of overall variance as any other  $k$ -dim. subspace.

- ▶  $K = 1$  (just a single Gaussian):

What is the  $k$ -dimensional PCA subspace?

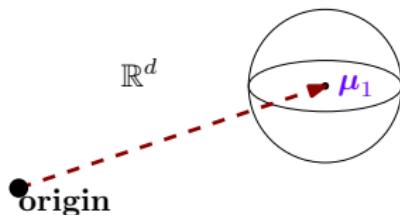


## Proof (continued)

**Key fact:**  $k$ -dimensional PCA subspace (based on  $\mathbb{E}(\mathbf{X} \otimes \mathbf{X})$ ) captures as much of overall variance as any other  $k$ -dim. subspace.

- ▶  $K = 1$  (just a single Gaussian):

What is the  $k$ -dimensional PCA subspace?



**Answer:** any  $k$ -dim. subspace containing  $\mu_1$ .

## Proof (continued)

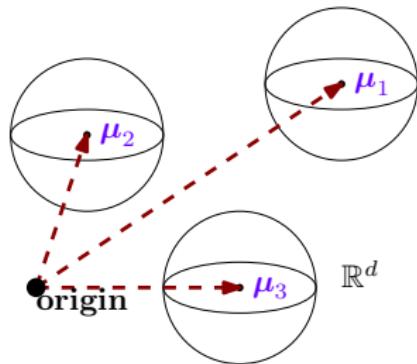
**Key fact:**  $k$ -dimensional PCA subspace (based on  $\mathbb{E}(\mathbf{X} \otimes \mathbf{X})$ ) captures as much of overall variance as any other  $k$ -dim. subspace.

- ▶ General  $K$  (mixture of  $K$  Gaussians):  
What is the  $K$ -dimensional PCA subspace?

## Proof (continued)

**Key fact:**  $k$ -dimensional PCA subspace (based on  $\mathbb{E}(\mathbf{X} \otimes \mathbf{X})$ ) captures as much of overall variance as any other  $k$ -dim. subspace.

- ▶ General  $K$  (mixture of  $K$  Gaussians):  
What is the  $K$ -dimensional PCA subspace?

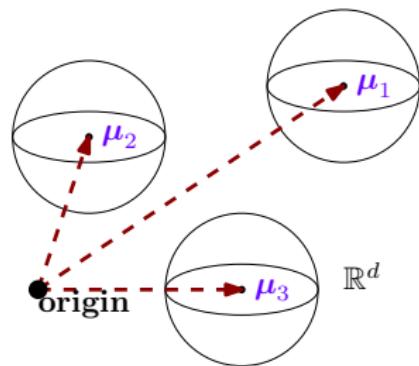


$$\mathbb{E}(\mathbf{X} \otimes \mathbf{X}) = \sum_{t=1}^K \pi_t \cdot \mu_t \otimes \mu_t + \sigma^2 \mathbf{I}_d .$$

## Proof (continued)

**Key fact:**  $k$ -dimensional PCA subspace (based on  $\mathbb{E}(\mathbf{X} \otimes \mathbf{X})$ ) captures as much of overall variance as any other  $k$ -dim. subspace.

- ▶ General  $K$  (mixture of  $K$  Gaussians):  
What is the  $K$ -dimensional PCA subspace?



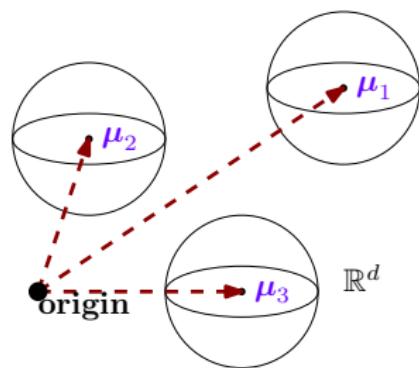
$$\mathbb{E}(\mathbf{X} \otimes \mathbf{X}) = \sum_{t=1}^K \pi_t \cdot \mu_t \otimes \mu_t + \sigma^2 \mathbf{I}_d .$$

**Answer:** any  $K$ -dim. subspace containing  $\mu_1, \dots, \mu_K$ . □

## Proof (continued)

**Key fact:**  $k$ -dimensional PCA subspace (based on  $\mathbb{E}(\mathbf{X} \otimes \mathbf{X})$ ) captures as much of overall variance as any other  $k$ -dim. subspace.

- ▶ General  $K$  (mixture of  $K$  Gaussians):  
What is the  $K$ -dimensional PCA subspace?



$$\mathbb{E}(\mathbf{X} \otimes \mathbf{X}) = \sum_{t=1}^K \pi_t \cdot \mu_t \otimes \mu_t + \sigma^2 \mathbf{I}_d.$$

**Answer:** any  $K$ -dim. subspace containing  $\mu_1, \dots, \mu_K$ . □

How does this help with learning mixtures of Gaussians?

## Use of moments for mixtures of spherical Gaussians

**Separation** (Dasgupta, 1999):

# standard deviations between component means

$$\text{sep} := \min_{i \neq j} \frac{\|\boldsymbol{\mu}_i - \boldsymbol{\mu}_j\|}{\sigma}.$$

## Use of moments for mixtures of spherical Gaussians

**Separation** (Dasgupta, 1999):

# standard deviations between component means

$$\text{sep} := \min_{i \neq j} \frac{\|\boldsymbol{\mu}_i - \boldsymbol{\mu}_j\|}{\sigma}.$$

► (Dasgupta & Schulman, 2000):

Distance-based clustering (e.g., EM) works when  $\text{sep} \gtrsim d^{1/4}$ .

# Use of moments for mixtures of spherical Gaussians

**Separation** (Dasgupta, 1999):

# standard deviations between component means

$$\text{sep} := \min_{i \neq j} \frac{\|\boldsymbol{\mu}_i - \boldsymbol{\mu}_j\|}{\sigma}.$$

► (Dasgupta & Schulman, 2000):

Distance-based clustering (e.g., EM) works when  $\text{sep} \gtrsim d^{1/4}$ .

► (Vempala & Wang, 2002):

Problem becomes  $K$ -dimensional via PCA (assume  $K \leq d$ ).

Required separation reduced to  $\text{sep} \gtrsim K^{1/4}$ .

# Use of moments for mixtures of spherical Gaussians

**Separation** (Dasgupta, 1999):

# standard deviations between component means

$$\text{sep} := \min_{i \neq j} \frac{\|\boldsymbol{\mu}_i - \boldsymbol{\mu}_j\|}{\sigma}.$$

► (Dasgupta & Schulman, 2000):

Distance-based clustering (e.g., EM) works when  $\text{sep} \gtrsim d^{1/4}$ .

► (Vempala & Wang, 2002):

Problem becomes  $K$ -dimensional via PCA (assume  $K \leq d$ ).

Required separation reduced to  $\text{sep} \gtrsim K^{1/4}$ .

---

**Third-order moments** identify the mixture distribution when

$\{\boldsymbol{\mu}_t\}_{t=1}^K$  are lin. indpt.;  $\text{sep}$  may be arbitrarily close to zero.

# Use of moments for mixtures of spherical Gaussians

**Separation** (Dasgupta, 1999):

# standard deviations between component means

$$\text{sep} := \min_{i \neq j} \frac{\|\boldsymbol{\mu}_i - \boldsymbol{\mu}_j\|}{\sigma}.$$

► (Dasgupta & Schulman, 2000):

Distance-based clustering (e.g., EM) works when  $\text{sep} \gtrsim d^{1/4}$ .

► (Vempala & Wang, 2002):

Problem becomes  $K$ -dimensional via PCA (assume  $K \leq d$ ).

Required separation reduced to  $\text{sep} \gtrsim K^{1/4}$ .

---

**Third-order moments** identify the mixture distribution when

$\{\boldsymbol{\mu}_t\}_{t=1}^K$  are lin. indpt.;  $\text{sep}$  may be arbitrarily close to zero.

(Belkin & Sinha, 2010; Moitra & Valiant, 2010):

General Gaussians & no minimum  $\text{sep}$ , but  $K$ th-order moments.

# Third-order moments of spherical Gaussian mixtures

**Generative process:**

$$\mathbf{X} = \mathbf{Y} + \sigma \mathbf{Z}$$

where  $\Pr(\mathbf{Y} = \boldsymbol{\mu}_t) = \pi_t$ , and  $\mathbf{Z} \sim \text{Normal}(\mathbf{0}, \mathbf{I}_d)$ ,  $\mathbf{Y} \perp\!\!\!\perp \mathbf{Z}$ .

---

Third-order moment tensor:

$$\mathbb{E}(\mathbf{X}^{\otimes 3}) = \mathbb{E}((\mathbf{Y} + \sigma \mathbf{Z})^{\otimes 3})$$

# Third-order moments of spherical Gaussian mixtures

**Generative process:**

$$\mathbf{X} = \mathbf{Y} + \sigma \mathbf{Z}$$

where  $\Pr(\mathbf{Y} = \boldsymbol{\mu}_t) = \pi_t$ , and  $\mathbf{Z} \sim \text{Normal}(\mathbf{0}, \mathbf{I}_d)$ ,  $\mathbf{Y} \perp\!\!\!\perp \mathbf{Z}$ .

---

Third-order moment tensor:

$$\begin{aligned}\mathbb{E}(\mathbf{X}^{\otimes 3}) &= \mathbb{E}((\mathbf{Y} + \sigma \mathbf{Z})^{\otimes 3}) \\ &= \mathbb{E}(\mathbf{Y}^{\otimes 3}) + \sigma^2 \mathbb{E}(\mathbf{Y} \otimes \mathbf{Z} \otimes \mathbf{Z} + \mathbf{Z} \otimes \mathbf{Y} \otimes \mathbf{Z} + \mathbf{Z} \otimes \mathbf{Z} \otimes \mathbf{Y})\end{aligned}$$

# Third-order moments of spherical Gaussian mixtures

**Generative process:**

$$\mathbf{X} = \mathbf{Y} + \sigma \mathbf{Z}$$

where  $\Pr(\mathbf{Y} = \boldsymbol{\mu}_t) = \pi_t$ , and  $\mathbf{Z} \sim \text{Normal}(\mathbf{0}, \mathbf{I}_d)$ ,  $\mathbf{Y} \perp\!\!\!\perp \mathbf{Z}$ .

---

Third-order moment tensor:

$$\begin{aligned}\mathbb{E}(\mathbf{X}^{\otimes 3}) &= \mathbb{E}((\mathbf{Y} + \sigma \mathbf{Z})^{\otimes 3}) \\ &= \mathbb{E}(\mathbf{Y}^{\otimes 3}) + \sigma^2 \mathbb{E}(\mathbf{Y} \otimes \mathbf{Z} \otimes \mathbf{Z} + \mathbf{Z} \otimes \mathbf{Y} \otimes \mathbf{Z} + \mathbf{Z} \otimes \mathbf{Z} \otimes \mathbf{Y}) \\ &= \sum_{t=1}^K \pi_t \cdot \boldsymbol{\mu}_t^{\otimes 3} + \sigma^2 \tau(\boldsymbol{\mu}).\end{aligned}$$

(Above,  $\boldsymbol{\mu} = \mathbb{E}(\mathbf{X})$  and  $\tau(\boldsymbol{\mu})$  is a third-order tensor involving only  $\boldsymbol{\mu}$ .)

# Third-order moments of spherical Gaussian mixtures

**Generative process:**

$$\mathbf{X} = \mathbf{Y} + \sigma \mathbf{Z}$$

where  $\Pr(\mathbf{Y} = \boldsymbol{\mu}_t) = \pi_t$ , and  $\mathbf{Z} \sim \text{Normal}(\mathbf{0}, \mathbf{I}_d)$ ,  $\mathbf{Y} \perp\!\!\!\perp \mathbf{Z}$ .

---

Third-order moment tensor:

$$\begin{aligned}\mathbb{E}(\mathbf{X}^{\otimes 3}) &= \mathbb{E}((\mathbf{Y} + \sigma \mathbf{Z})^{\otimes 3}) \\ &= \mathbb{E}(\mathbf{Y}^{\otimes 3}) + \sigma^2 \mathbb{E}(\mathbf{Y} \otimes \mathbf{Z} \otimes \mathbf{Z} + \mathbf{Z} \otimes \mathbf{Y} \otimes \mathbf{Z} + \mathbf{Z} \otimes \mathbf{Z} \otimes \mathbf{Y}) \\ &= \sum_{t=1}^K \pi_t \cdot \boldsymbol{\mu}_t^{\otimes 3} + \sigma^2 \tau(\boldsymbol{\mu}).\end{aligned}$$

(Above,  $\boldsymbol{\mu} = \mathbb{E}(\mathbf{X})$  and  $\tau(\boldsymbol{\mu})$  is a third-order tensor involving only  $\boldsymbol{\mu}$ .)

**Exercise:** find explicit formula for  $\tau(\boldsymbol{\mu})$ .

# Tensor decomposition for spherical Gaussian mixtures

(H. & Kakade, 2013)

$$\mathbb{E} \left( \mathbf{X}^{\otimes 3} \right) = \sum_{t=1}^K \pi_t \cdot \boldsymbol{\mu}_t^{\otimes 3} + \sigma^2 \tau(\boldsymbol{\mu}).$$

# Tensor decomposition for spherical Gaussian mixtures

(H. & Kakade, 2013)

$$\mathbb{E}(\mathbf{X}^{\otimes 3}) = \sum_{t=1}^K \pi_t \cdot \boldsymbol{\mu}_t^{\otimes 3} + \sigma^2 \tau(\boldsymbol{\mu}).$$

**Claim:**  $\boldsymbol{\mu}$  &  $\sigma^2$  are simple functions of  $\mathbb{E}(\mathbf{X})$  &  $\mathbb{E}(\mathbf{X} \otimes \mathbf{X})$ .

# Tensor decomposition for spherical Gaussian mixtures

(H. & Kakade, 2013)

$$\mathbb{E}(\mathbf{X}^{\otimes 3}) = \sum_{t=1}^K \pi_t \cdot \boldsymbol{\mu}_t^{\otimes 3} + \sigma^2 \tau(\boldsymbol{\mu}).$$

**Claim:**  $\boldsymbol{\mu}$  &  $\sigma^2$  are simple functions of  $\mathbb{E}(\mathbf{X})$  &  $\mathbb{E}(\mathbf{X} \otimes \mathbf{X})$ .

**Claim:** If  $\{\boldsymbol{\mu}_t\}_{t=1}^K$  are linearly independent and all  $\pi_t > 0$ , then  $\{(\boldsymbol{\mu}_t, \pi_t)\}_{t=1}^K$  are identifiable from

$$\mathbf{T} := \mathbb{E}(\mathbf{X}^{\otimes 3}) - \sigma^2 \tau(\boldsymbol{\mu}) = \sum_{t=1}^K \pi_t \cdot \boldsymbol{\mu}_t^{\otimes 3}.$$

# Tensor decomposition for spherical Gaussian mixtures

(H. & Kakade, 2013)

$$\mathbb{E}(\mathbf{X}^{\otimes 3}) = \sum_{t=1}^K \pi_t \cdot \boldsymbol{\mu}_t^{\otimes 3} + \sigma^2 \tau(\boldsymbol{\mu}).$$

**Claim:**  $\boldsymbol{\mu}$  &  $\sigma^2$  are simple functions of  $\mathbb{E}(\mathbf{X})$  &  $\mathbb{E}(\mathbf{X} \otimes \mathbf{X})$ .

**Claim:** If  $\{\boldsymbol{\mu}_t\}_{t=1}^K$  are linearly independent and all  $\pi_t > 0$ , then  $\{(\boldsymbol{\mu}_t, \pi_t)\}_{t=1}^K$  are identifiable from

$$\mathbf{T} := \mathbb{E}(\mathbf{X}^{\otimes 3}) - \sigma^2 \tau(\boldsymbol{\mu}) = \sum_{t=1}^K \pi_t \cdot \boldsymbol{\mu}_t^{\otimes 3}.$$

Can use tensor decomposition to recover  $\{(\boldsymbol{\mu}_t, \pi_t)\}_{t=1}^K$  from  $\mathbf{T}$ .

## Even more Gaussian mixtures

**Note:** Linear independence condition on  $\{\mu_t\}_{t=1}^K$  requires  $K \leq d$ .

## Even more Gaussian mixtures

**Note:** Linear independence condition on  $\{\mu_t\}_{t=1}^K$  requires  $K \leq d$ .

- ▶ (Anderson, Belkin, Goyal, Rademacher, & Voss, 2014),  
(Bhaskara, Charikar, Moitra, & Vijayaraghavan, 2014)

Mixtures of  $d^{O(1)}$  Gaussians (w/ simple or known covariance)  
via **smoothed analysis** and  $O(1)$ -order moments.

## Even more Gaussian mixtures

**Note:** Linear independence condition on  $\{\mu_t\}_{t=1}^K$  requires  $K \leq d$ .

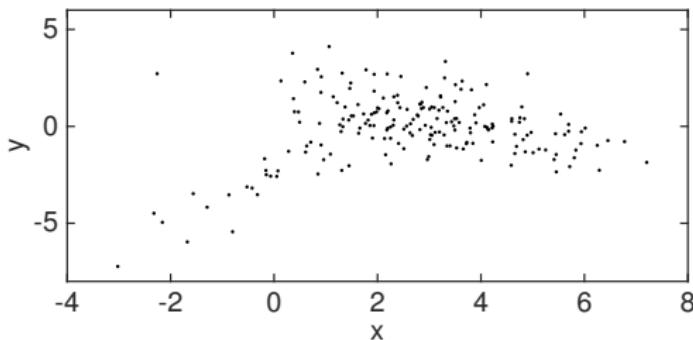
- ▶ (Anderson, Belkin, Goyal, Rademacher, & Voss, 2014),  
(Bhaskara, Charikar, Moitra, & Vijayaraghavan, 2014)  
Mixtures of  $d^{O(1)}$  Gaussians (w/ simple or known covariance)  
via **smoothed analysis** and  $O(1)$ -order moments.
- ▶ (Ge, Huang, & Kakade, 2015)  
Also with **unknown covariances of arbitrary shape**.

## Mixture model #2: Mixtures of linear regressions

$$\begin{aligned} H &\sim \text{Categorical}(\pi_1, \pi_2, \dots, \pi_K) \quad (\text{hidden}) ; \\ \boldsymbol{X} &\sim \text{Normal}(\boldsymbol{\mu}, \boldsymbol{\Sigma}) ; \\ Y \mid H = t, \boldsymbol{X} = \boldsymbol{x} &\sim \text{Normal}(\langle \boldsymbol{\beta}_t, \boldsymbol{x} \rangle, \sigma^2) . \end{aligned}$$

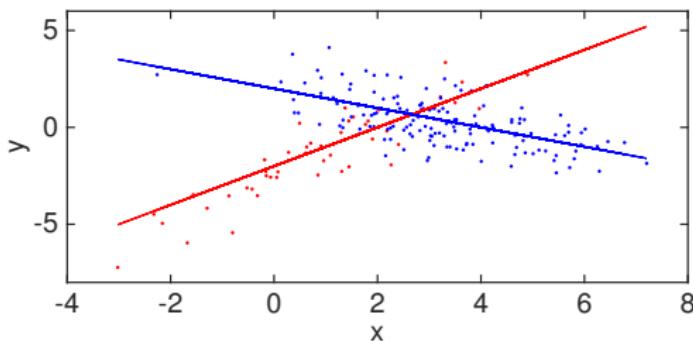
## Mixture model #2: Mixtures of linear regressions

$$\begin{aligned} H &\sim \text{Categorical}(\pi_1, \pi_2, \dots, \pi_K) \quad (\text{hidden}); \\ \boldsymbol{X} &\sim \text{Normal}(\boldsymbol{\mu}, \boldsymbol{\Sigma}); \\ Y \mid H = t, \boldsymbol{X} = \boldsymbol{x} &\sim \text{Normal}(\langle \boldsymbol{\beta}_t, \boldsymbol{x} \rangle, \sigma^2). \end{aligned}$$



## Mixture model #2: Mixtures of linear regressions

$$\begin{aligned} H &\sim \text{Categorical}(\pi_1, \pi_2, \dots, \pi_K) \quad (\text{hidden}); \\ \boldsymbol{X} &\sim \text{Normal}(\boldsymbol{\mu}, \boldsymbol{\Sigma}); \\ Y \mid H = t, \boldsymbol{X} = \boldsymbol{x} &\sim \text{Normal}(\langle \boldsymbol{\beta}_t, \boldsymbol{x} \rangle, \sigma^2). \end{aligned}$$



## Use of moments for mixtures of linear regressions

**Second-order moments** (assume  $\mathbf{X} \sim \text{Normal}(\mathbf{0}, \mathbf{I}_d)$ ):

$$\mathbb{E}(\mathbf{Y}^2 \mathbf{X} \mathbf{X}^\top) = 2 \sum_{t=1}^K \pi_t \cdot \boldsymbol{\beta}_t \boldsymbol{\beta}_t^\top + \left( \sigma^2 + \sum_{t=1}^K \pi_t \cdot \|\boldsymbol{\beta}_t\|^2 \right) \mathbf{I}_d.$$

## Use of moments for mixtures of linear regressions

**Second-order moments** (assume  $\mathbf{X} \sim \text{Normal}(\mathbf{0}, \mathbf{I}_d)$ ):

$$\mathbb{E}(\mathbf{Y}^2 \mathbf{X} \mathbf{X}^\top) = 2 \sum_{t=1}^K \pi_t \cdot \boldsymbol{\beta}_t \boldsymbol{\beta}_t^\top + \left( \sigma^2 + \sum_{t=1}^K \pi_t \cdot \|\boldsymbol{\beta}_t\|^2 \right) \mathbf{I}_d.$$

- ▶ Span of top  $K$  eigenvectors of  $\mathbb{E}(\mathbf{Y}^2 \mathbf{X} \mathbf{X}^\top)$  contains  $\{\boldsymbol{\beta}_t\}_{t=1}^K$ .

## Use of moments for mixtures of linear regressions

**Second-order moments** (assume  $\mathbf{X} \sim \text{Normal}(\mathbf{0}, \mathbf{I}_d)$ ):

$$\mathbb{E}(\mathbf{Y}^2 \mathbf{X} \mathbf{X}^\top) = 2 \sum_{t=1}^K \pi_t \cdot \boldsymbol{\beta}_t \boldsymbol{\beta}_t^\top + \left( \sigma^2 + \sum_{t=1}^K \pi_t \cdot \|\boldsymbol{\beta}_t\|^2 \right) \mathbf{I}_d.$$

- ▶ Span of top  $K$  eigenvectors of  $\mathbb{E}(\mathbf{Y}^2 \mathbf{X} \mathbf{X}^\top)$  contains  $\{\boldsymbol{\beta}_t\}_{t=1}^K$ .
- ▶ Using Stein's identity (1973), similar approach works for GLMs (Sun, Ioannidis, & Montanari, 2013).

# Use of moments for mixtures of linear regressions

**Second-order moments** (assume  $\mathbf{X} \sim \text{Normal}(\mathbf{0}, \mathbf{I}_d)$ ):

$$\mathbb{E}(\mathbf{Y}^2 \mathbf{X} \mathbf{X}^\top) = 2 \sum_{t=1}^K \pi_t \cdot \boldsymbol{\beta}_t \boldsymbol{\beta}_t^\top + \left( \sigma^2 + \sum_{t=1}^K \pi_t \cdot \|\boldsymbol{\beta}_t\|^2 \right) \mathbf{I}_d.$$

- ▶ Span of top  $K$  eigenvectors of  $\mathbb{E}(\mathbf{Y}^2 \mathbf{X} \mathbf{X}^\top)$  contains  $\{\boldsymbol{\beta}_t\}_{t=1}^K$ .
- ▶ Using Stein's identity (1973), similar approach works for GLMs (Sun, Ioannidis, & Montanari, 2013).

**Tensor decomposition approach:**

Can recover parameters  $\{(\boldsymbol{\beta}_t, \pi_t)\}_{t=1}^K$  with higher-order moments (Chaganty & Liang, 2013; Yi, Caramanis, & Sanghavi, 2014, 2016).

# Use of moments for mixtures of linear regressions

**Second-order moments** (assume  $\mathbf{X} \sim \text{Normal}(\mathbf{0}, \mathbf{I}_d)$ ):

$$\mathbb{E}(\mathbf{Y}^2 \mathbf{X} \mathbf{X}^\top) = 2 \sum_{t=1}^K \pi_t \cdot \boldsymbol{\beta}_t \boldsymbol{\beta}_t^\top + \left( \sigma^2 + \sum_{t=1}^K \pi_t \cdot \|\boldsymbol{\beta}_t\|^2 \right) \mathbf{I}_d.$$

- ▶ Span of top  $K$  eigenvectors of  $\mathbb{E}(\mathbf{Y}^2 \mathbf{X} \mathbf{X}^\top)$  contains  $\{\boldsymbol{\beta}_t\}_{t=1}^K$ .
- ▶ Using Stein's identity (1973), similar approach works for GLMs (Sun, Ioannidis, & Montanari, 2013).

**Tensor decomposition approach:**

Can recover parameters  $\{(\boldsymbol{\beta}_t, \pi_t)\}_{t=1}^K$  with higher-order moments (Chaganty & Liang, 2013; Yi, Caramanis, & Sanghavi, 2014, 2016).

Also for GLMs, via Stein's identity (Sedghi & Anandkumar, 2014).

## Recap: mixtures of Gaussians and linear regressions

- ▶ Parameters of Gaussian mixture models and related models (satisfying linear independence condition) can be efficiently recovered from  $O(1)$ -order moments.

## Recap: mixtures of Gaussians and linear regressions

- ▶ Parameters of Gaussian mixture models and related models (satisfying linear independence condition) can be efficiently recovered from  $O(1)$ -order moments.
- ▶ Exploit distributional properties to determine usable moments.

## Recap: mixtures of Gaussians and linear regressions

- ▶ Parameters of Gaussian mixture models and related models (satisfying linear independence condition) can be efficiently recovered from  $O(1)$ -order moments.
- ▶ Exploit distributional properties to determine usable moments.
- ▶ *Smoothed analysis*: avoid linear independence condition for “most” mixture distributions.

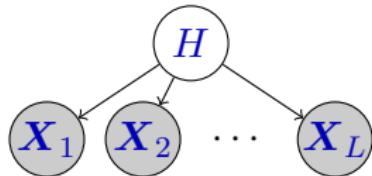
## Recap: mixtures of Gaussians and linear regressions

- ▶ Parameters of Gaussian mixture models and related models (satisfying linear independence condition) can be efficiently recovered from  $O(1)$ -order moments.
- ▶ Exploit distributional properties to determine usable moments.
- ▶ *Smoothed analysis*: avoid linear independence condition for “most” mixture distributions.

**Next:** Multi-view approach to finding usable moments.

# Multi-view interpretation of topic model

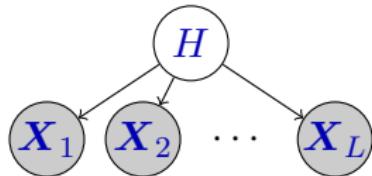
**Recall:** Topic model for single-topic documents



$K$  topics (dists. over words)  $\{\mathbf{P}_t\}_{t=1}^K$ .  
Pick topic  $H = t$  with prob.  $w_t$  (hidden).  
Word tokens  $X_1, X_2, \dots, X_L \stackrel{\text{iid}}{\sim} \mathbf{P}_H$ .

# Multi-view interpretation of topic model

**Recall:** Topic model for single-topic documents



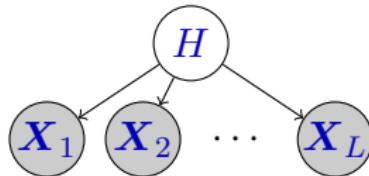
$K$  topics (dists. over words)  $\{\mathbf{P}_t\}_{t=1}^K$ .  
Pick topic  $H = t$  with prob.  $w_t$  (hidden).  
Word tokens  $\mathbf{X}_1, \mathbf{X}_2, \dots, \mathbf{X}_L \stackrel{\text{iid}}{\sim} \mathbf{P}_H$ .

**Key property:**

$\mathbf{X}_1, \mathbf{X}_2, \dots, \mathbf{X}_L$  conditionally independent given  $H$ .

# Multi-view interpretation of topic model

**Recall:** Topic model for single-topic documents



$K$  topics (dists. over words)  $\{\mathbf{P}_t\}_{t=1}^K$ .  
Pick topic  $H = t$  with prob.  $w_t$  (hidden).  
Word tokens  $\mathbf{X}_1, \mathbf{X}_2, \dots, \mathbf{X}_L \stackrel{\text{iid}}{\sim} \mathbf{P}_H$ .

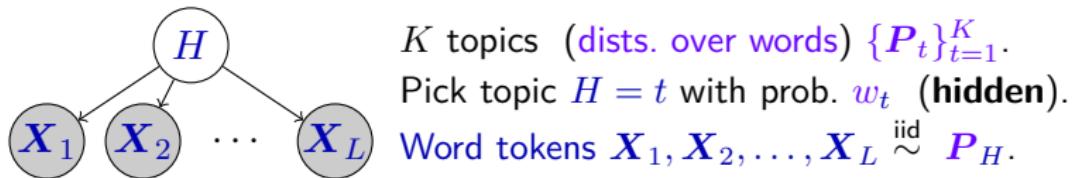
**Key property:**

$\mathbf{X}_1, \mathbf{X}_2, \dots, \mathbf{X}_L$  conditionally independent given  $H$ .

Each word token  $\mathbf{X}_i$  provides new “view” of hidden variable  $H$ .

# Multi-view interpretation of topic model

**Recall:** Topic model for single-topic documents



**Key property:**

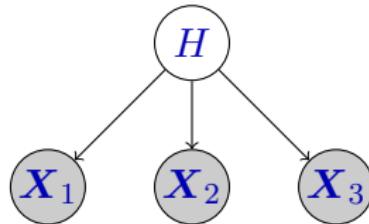
$\mathbf{X}_1, \mathbf{X}_2, \dots, \mathbf{X}_L$  conditionally independent given  $H$ .

Each word token  $\mathbf{X}_i$  provides new “view” of hidden variable  $H$ .

**Some previous analyses:**

- ▶ (Blum & Mitchell, 1998)  
*Co-training* in semi-supervised learning.
- ▶ (Chaudhuri, Kakade, Livescu, & Sridharan, 2009)  
Multi-view Gaussian mixture models.

# Multi-view mixture model



View 1:  $X_1$

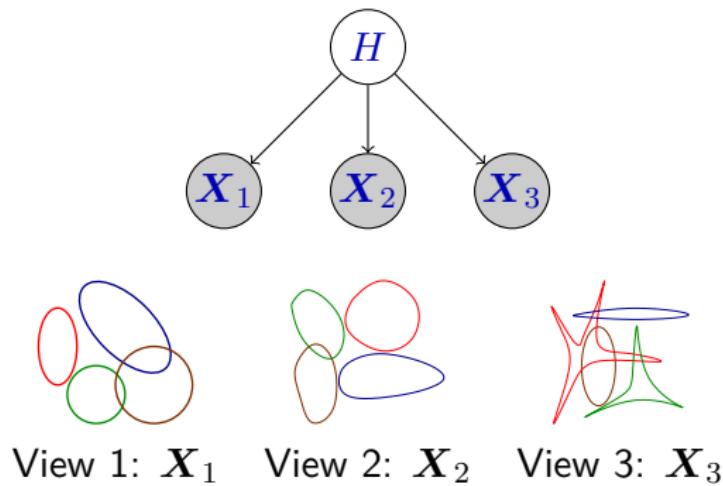


View 2:  $X_2$

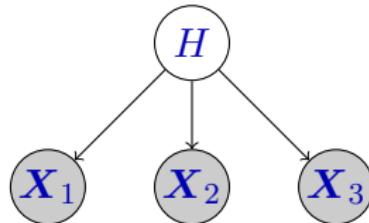


View 3:  $X_3$

## Multi-view mixture model



## Multi-view mixture model

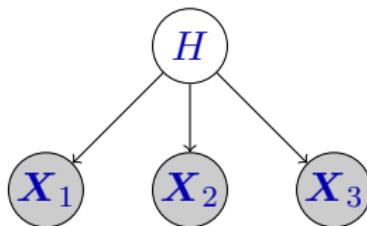


$$\mathbb{E}(\mathbf{X}_1 \otimes \mathbf{X}_2 \otimes \mathbf{X}_3) = \sum_{t=1}^K \pi_t \cdot \boldsymbol{\mu}_t^{(1)} \otimes \boldsymbol{\mu}_t^{(2)} \otimes \boldsymbol{\mu}_t^{(3)}$$

where  $\boldsymbol{\mu}_t^{(i)} = \mathbb{E}[\mathbf{X}_i \mid H = t]$ ,

$\pi_t = \Pr(H = t)$ .

## Multi-view mixture model



$$\mathbb{E}(\mathbf{X}_1 \otimes \mathbf{X}_2 \otimes \mathbf{X}_3) = \sum_{t=1}^K \pi_t \cdot \boldsymbol{\mu}_t^{(1)} \otimes \boldsymbol{\mu}_t^{(2)} \otimes \boldsymbol{\mu}_t^{(3)}$$

$$\text{where } \boldsymbol{\mu}_t^{(i)} = \mathbb{E}[\mathbf{X}_i \mid H = t],$$

$$\pi_t = \Pr(H = t).$$

**Tensor decomposition approach** works in this asymmetric case as long as  $\{\boldsymbol{\mu}_t^{(j)}\}_{t=1}^K$  are lin. indpt. for each  $j$ , and all  $\pi_t > 0$ .

## Examples of multi-view mixture models

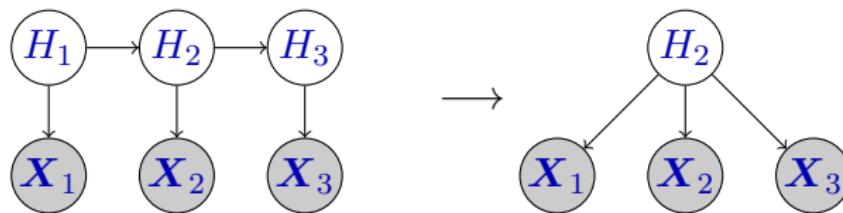
(Mossel & Roch, 2006; Anandkumar, H., & Kakade, 2012)

1. Mixtures of high-dimensional product distributions.  
(E.g., mixtures of axis-aligned Gaussians, other topic models.)

## Examples of multi-view mixture models

(Mossel & Roch, 2006; Anandkumar, H., & Kakade, 2012)

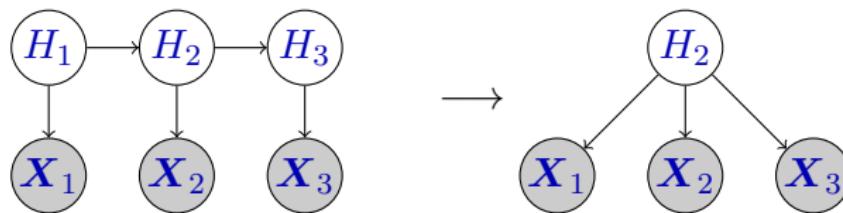
1. Mixtures of high-dimensional product distributions.  
(E.g., mixtures of axis-aligned Gaussians, other topic models.)
2. Hidden Markov models.



## Examples of multi-view mixture models

(Mossel & Roch, 2006; Anandkumar, H., & Kakade, 2012)

1. Mixtures of high-dimensional product distributions.  
(E.g., mixtures of axis-aligned Gaussians, other topic models.)
2. Hidden Markov models.



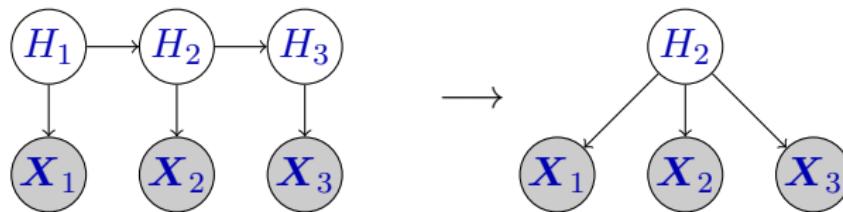
3. Phylogenetic trees.

- ▶  $X_1, X_2, X_3$ : genes of three extant species.
- ▶  $H$ : LCA of most closely related pair of species.

## Examples of multi-view mixture models

(Mossel & Roch, 2006; Anandkumar, H., & Kakade, 2012)

1. Mixtures of high-dimensional product distributions.  
(E.g., mixtures of axis-aligned Gaussians, other topic models.)
2. Hidden Markov models.

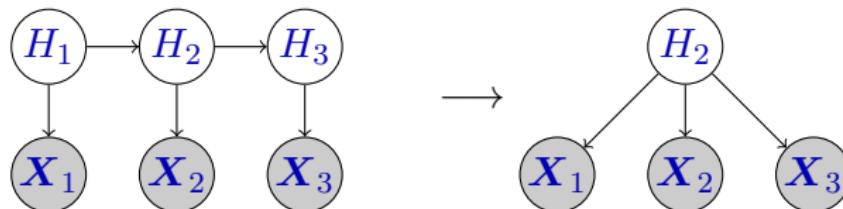


3. Phylogenetic trees.
  - ▶  $X_1, X_2, X_3$ : genes of three extant species.
  - ▶  $H$ : LCA of most closely related pair of species.
4. ...

# Examples of multi-view mixture models

(Mossel & Roch, 2006; Anandkumar, H., & Kakade, 2012)

1. Mixtures of high-dimensional product distributions.  
(E.g., mixtures of axis-aligned Gaussians, other topic models.)
2. Hidden Markov models.



3. Phylogenetic trees.
  - ▶  $X_1, X_2, X_3$ : genes of three extant species.
  - ▶  $H$ : LCA of most closely related pair of species.
4. ...

**Next:** Single index models.

## Single-index models

$$\begin{aligned}\mathbf{X} &\sim \text{Normal}(\mathbf{0}, \mathbf{I}) ; \\ Y \mid \mathbf{X} = \mathbf{x} &\sim \text{Normal}(g(\langle \boldsymbol{\beta}, \mathbf{x} \rangle), \sigma^2) .\end{aligned}$$

Here,  $g: \mathbb{R} \rightarrow \mathbb{R}$  is the *link function*.

# Single-index models

$$\mathbf{X} \sim \text{Normal}(\mathbf{0}, \mathbf{I});$$

$$Y \mid \mathbf{X} = \mathbf{x} \sim \text{Normal}(g(\langle \boldsymbol{\beta}, \mathbf{x} \rangle), \sigma^2).$$

Here,  $g: \mathbb{R} \rightarrow \mathbb{R}$  is the *link function*.

- ▶ **Phase retrieval** (real signals): assume  $g(z) = z^2$ .
- ▶ **1-bit compressed sensing**: assume  $g(z) = \text{sign}(z)$ .
- ▶ **Isotonic regression**: assume  $g$  is monotone (e.g.,  $g' \geq 0$ ).
- ▶ **Convex regression**: assume  $g$  is convex (e.g.,  $g'' \geq 0$ ).
- ▶ ...

# Single-index models

$$\mathbf{X} \sim \text{Normal}(\mathbf{0}, \mathbf{I});$$

$$Y \mid \mathbf{X} = \mathbf{x} \sim \text{Normal}(g(\langle \boldsymbol{\beta}, \mathbf{x} \rangle), \sigma^2).$$

Here,  $g: \mathbb{R} \rightarrow \mathbb{R}$  is the *link function*.

- ▶ **Phase retrieval** (real signals): assume  $g(z) = z^2$ .
- ▶ **1-bit compressed sensing**: assume  $g(z) = \text{sign}(z)$ .
- ▶ **Isotonic regression**: assume  $g$  is monotone (e.g.,  $g' \geq 0$ ).
- ▶ **Convex regression**: assume  $g$  is convex (e.g.,  $g'' \geq 0$ ).
- ▶ ...

When  $g$  is unknown, model is generally called **single-index model**.

# Single-index models

$$\mathbf{X} \sim \text{Normal}(\mathbf{0}, \mathbf{I});$$

$$Y \mid \mathbf{X} = \mathbf{x} \sim \text{Normal}(g(\langle \boldsymbol{\beta}, \mathbf{x} \rangle), \sigma^2).$$

Here,  $g: \mathbb{R} \rightarrow \mathbb{R}$  is the *link function*.

- ▶ **Phase retrieval** (real signals): assume  $g(z) = z^2$ .
- ▶ **1-bit compressed sensing**: assume  $g(z) = \text{sign}(z)$ .
- ▶ **Isotonic regression**: assume  $g$  is monotone (e.g.,  $g' \geq 0$ ).
- ▶ **Convex regression**: assume  $g$  is convex (e.g.,  $g'' \geq 0$ ).
- ▶ ...

When  $g$  is unknown, model is generally called **single-index model**.

**Semi-parametric estimation**: regard  $g$  as nuisance parameter;  
focus on estimating  $\boldsymbol{\beta}$ .

## Aside: symmetric tensors and homogeneous polynomials

Recall formula for tensor function value:

$$\mathbf{T}(\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(p)}) = \sum_{i_1, \dots, i_p} \mathbf{T}_{i_1, \dots, i_p} \cdot x_{i_1}^{(1)} \cdots x_{i_p}^{(p)}.$$

## Aside: symmetric tensors and homogeneous polynomials

Recall formula for tensor function value:

$$\mathbf{T}(\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(p)}) = \sum_{i_1, \dots, i_p} \mathbf{T}_{i_1, \dots, i_p} \cdot x_{i_1}^{(1)} \cdots x_{i_p}^{(p)}.$$

If  $\mathbf{T}$  is symmetric (i.e.,  $\mathbf{T}_{i_1, \dots, i_p} = \mathbf{T}_{\pi(i_1), \dots, \pi(i_p)}$  for any permutation  $\pi$ ), then evaluating at  $\mathbf{x}^{(1)} = \dots = \mathbf{x}^{(p)} = \mathbf{x}$  gives

$$\mathbf{T}(\mathbf{x}, \dots, \mathbf{x}) = p! \sum_{i_1 < \dots < i_p} \mathbf{T}_{i_1, \dots, i_p} \cdot x_{i_1} \cdots x_{i_p},$$

which is just the formula for a degree- $p$  homogeneous polynomial.

## Aside: symmetric tensors and homogeneous polynomials

Recall formula for tensor function value:

$$\mathbf{T}(\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(p)}) = \sum_{i_1, \dots, i_p} \mathbf{T}_{i_1, \dots, i_p} \cdot x_{i_1}^{(1)} \cdots x_{i_p}^{(p)}.$$

If  $\mathbf{T}$  is symmetric (i.e.,  $\mathbf{T}_{i_1, \dots, i_p} = \mathbf{T}_{\pi(i_1), \dots, \pi(i_p)}$  for any permutation  $\pi$ ), then evaluating at  $\mathbf{x}^{(1)} = \dots = \mathbf{x}^{(p)} = \mathbf{x}$  gives

$$\mathbf{T}(\mathbf{x}, \dots, \mathbf{x}) = p! \sum_{i_1 < \dots < i_p} \mathbf{T}_{i_1, \dots, i_p} \cdot x_{i_1} \cdots x_{i_p},$$

which is just the formula for a degree- $p$  homogeneous polynomial.

$p$ -th order symmetric tensors  $\simeq$  degree- $p$  homogeneous polynomials.

## Using orthogonal polynomials

(Dudeja & H., 2018)

Let  $H_p: \mathbb{R} \rightarrow \mathbb{R}$  denote the degree- $p$  Hermite polynomial.

Assume (for  $Z \sim \text{Normal}(0, 1)$ ):

- ▶  $\mathbb{E}[g(Z)^2] = 1$  (normalization—this is WLOG);
- ▶  $\mathbb{E}[g'(Z)^2] \geq \epsilon$  (necessary for identifiability);
- ▶  $g$  is smooth and  $\mathbb{E}[g''(Z)^2] = O(1)$ .

# Using orthogonal polynomials

(Dudeja & H., 2018)

Let  $H_p: \mathbb{R} \rightarrow \mathbb{R}$  denote the degree- $p$  Hermite polynomial.

Assume (for  $Z \sim \text{Normal}(0, 1)$ ):

- ▶  $\mathbb{E}[g(Z)^2] = 1$  (normalization—this is WLOG);
- ▶  $\mathbb{E}[g'(Z)^2] \geq \epsilon$  (necessary for identifiability);
- ▶  $g$  is smooth and  $\mathbb{E}[g''(Z)^2] = O(1)$ .

There exists  $p = O(1/\epsilon)$  such that

$$\mathbb{E}[\mathbf{Y} H_p(\langle \mathbf{v}, \mathbf{X} \rangle)] = (\lambda \boldsymbol{\beta}^{\otimes p})(\mathbf{v}), \quad \mathbf{v} \in \mathbb{R}^d$$

for some  $\lambda \neq 0$  with  $|\lambda| = \Omega(\epsilon/\sqrt{p})$ .

# Using orthogonal polynomials

(Dudeja & H., 2018)

Let  $H_p: \mathbb{R} \rightarrow \mathbb{R}$  denote the degree- $p$  Hermite polynomial.

Assume (for  $Z \sim \text{Normal}(0, 1)$ ):

- ▶  $\mathbb{E}[g(Z)^2] = 1$  (normalization—this is WLOG);
- ▶  $\mathbb{E}[g'(Z)^2] \geq \epsilon$  (necessary for identifiability);
- ▶  $g$  is smooth and  $\mathbb{E}[g''(Z)^2] = O(1)$ .

There exists  $p = O(1/\epsilon)$  such that

$$\mathbb{E}[Y H_p(\langle \mathbf{v}, \mathbf{X} \rangle)] = (\lambda \beta^{\otimes p})(\mathbf{v}), \quad \mathbf{v} \in \mathbb{R}^d$$

for some  $\lambda \neq 0$  with  $|\lambda| = \Omega(\epsilon/\sqrt{p})$ .

⇒ Get efficient algorithms for semi-parametric estimation of single-index model parameters, for very general link functions.

## Recap

- ▶ Parameters of many latent variable models (satisfying non-degeneracy conditions) can be efficiently recovered from  $O(1)$ -order moments.

## Recap

- ▶ Parameters of many latent variable models (satisfying non-degeneracy conditions) can be efficiently recovered from  $O(1)$ -order moments.
- ▶ Exploit distributional properties, multi-view structure, and other structure to determine usable moments.

## Recap

- ▶ Parameters of many latent variable models (satisfying non-degeneracy conditions) can be efficiently recovered from  $O(1)$ -order moments.
- ▶ Exploit distributional properties, multi-view structure, and other structure to determine usable moments.
- ▶ Estimation via **method-of-moments**:
  1. Estimate moments  $\rightarrow$  empirical moment tensor  $\hat{\mathbf{T}}$ .
  2. Approximately decompose  $\hat{\mathbf{T}}$   $\rightarrow$  parameter estimate  $\hat{\boldsymbol{\theta}}$ .

### 3. Error analysis (if time)

## Moment estimates

Estimation of  $\mathbb{E}[\mathbf{X}^{\otimes 3}]$  (say) from iid sample  $\{\mathbf{x}_i\}_{i=1}^n$ :

$$\widehat{\mathbb{E}}[\mathbf{X}^{\otimes 3}] := \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i^{\otimes 3}.$$

## Moment estimates

Estimation of  $\mathbb{E}[\mathbf{X}^{\otimes 3}]$  (say) from iid sample  $\{\mathbf{x}_i\}_{i=1}^n$ :

$$\widehat{\mathbb{E}}[\mathbf{X}^{\otimes 3}] := \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i^{\otimes 3}.$$

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

$$\|\mathbf{T}\| := \sup_{\|\mathbf{x}\|=\|\mathbf{y}\|=\|\mathbf{z}\|=1} \mathbf{T}(\mathbf{x}, \mathbf{y}, \mathbf{z}) \quad (\text{injective/“spectral” norm}),$$

$$\|\mathbf{T}\|_F := \left( \sum_{i,j,k} T_{i,j,k}^2 \right)^{1/2} \quad (\text{Frobenius norm}).$$

## Nearly orthogonally decomposable tensor

(Mu, H., & Goldfarb, 2015)

Let  $\varepsilon = \|\textcolor{red}{E}\|$  for  $\textcolor{red}{E} := \hat{\mathbf{T}} - \mathbf{T}$ .

**Claim:** Let  $\hat{\mathbf{v}} := \arg \max_{\|\mathbf{x}\|=1} \hat{\mathbf{T}}(\mathbf{x}, \mathbf{x}, \mathbf{x})$  and  $\hat{\lambda} := \hat{\mathbf{T}}(\hat{\mathbf{v}}, \hat{\mathbf{v}}, \hat{\mathbf{v}})$ .

Then

$$|\hat{\lambda} - \lambda_t| \leq \varepsilon, \quad \|\hat{\mathbf{v}} - \mathbf{v}_t\| \leq O\left(\frac{\varepsilon}{\lambda_t} + \left(\frac{\varepsilon}{\lambda_t}\right)^2\right)$$

for some  $t \in [d]$  with  $\lambda_t \geq \max_{t'} \lambda_{t'} - 2\varepsilon$ .

## Nearly orthogonally decomposable tensor

(Mu, H., & Goldfarb, 2015)

Let  $\varepsilon = \|\mathbf{E}\|$  for  $\mathbf{E} := \hat{\mathbf{T}} - \mathbf{T}$ .

**Claim:** Let  $\hat{\mathbf{v}} := \arg \max_{\|\mathbf{x}\|=1} \hat{\mathbf{T}}(\mathbf{x}, \mathbf{x}, \mathbf{x})$  and  $\hat{\lambda} := \hat{\mathbf{T}}(\hat{\mathbf{v}}, \hat{\mathbf{v}}, \hat{\mathbf{v}})$ .

Then

$$|\hat{\lambda} - \lambda_t| \leq \varepsilon, \quad \|\hat{\mathbf{v}} - \mathbf{v}_t\| \leq O\left(\frac{\varepsilon}{\lambda_t} + \left(\frac{\varepsilon}{\lambda_t}\right)^2\right)$$

for some  $t \in [d]$  with  $\lambda_t \geq \max_{t'} \lambda_{t'} - 2\varepsilon$ .

Many efficient algorithms for solving this approximately, when  $\varepsilon$  is small enough, like  $1/d$  or  $1/\sqrt{d}$  (e.g., Anandkumar, Ge, H., Kakade, & Telgarsky, 2014; Ma, Shi, & Steurer, 2016).

## Recall: greedy decomposition

(Zhang & Golub, 2001)

**Matching moments:**

$$\{(\hat{\mathbf{v}}_t, \hat{\lambda}_t)\}_{t=1}^d := \arg \min_{\{(\mathbf{x}_t, \sigma_t)\}_{t=1}^d} \left\| \mathbf{T} - \sum_{t=1}^d \sigma_t \cdot \mathbf{x}_t \otimes \mathbf{x}_t \otimes \mathbf{x}_t \right\|_F^2.$$

## Recall: greedy decomposition

(Zhang & Golub, 2001)

**Matching moments:**

$$\{(\hat{\mathbf{v}}_t, \hat{\lambda}_t)\}_{t=1}^d := \arg \min_{\{(\mathbf{x}_t, \sigma_t)\}_{t=1}^d} \left\| \mathbf{T} - \sum_{t=1}^d \sigma_t \cdot \mathbf{x}_t \otimes \mathbf{x}_t \otimes \mathbf{x}_t \right\|_F^2.$$

- ▶ Greedy approach:
  - ▶ Find best rank-1 approximation:

$$(\hat{\mathbf{v}}, \hat{\lambda}) := \arg \min_{\|\mathbf{x}\|=1, \sigma \geq 0} \|\mathbf{T} - \sigma \cdot \mathbf{x} \otimes \mathbf{x} \otimes \mathbf{x}\|_F^2.$$

- ▶ “Deflate”  $\mathbf{T} := \mathbf{T} - \hat{\lambda} \cdot \hat{\mathbf{v}} \otimes \hat{\mathbf{v}} \otimes \hat{\mathbf{v}}$  and repeat.

## Recall: greedy decomposition

(Zhang & Golub, 2001)

**Matching moments:**

$$\{(\hat{\mathbf{v}}_t, \hat{\lambda}_t)\}_{t=1}^d := \arg \min_{\{(\mathbf{x}_t, \sigma_t)\}_{t=1}^d} \left\| \mathbf{T} - \sum_{t=1}^d \sigma_t \cdot \mathbf{x}_t \otimes \mathbf{x}_t \otimes \mathbf{x}_t \right\|_F^2.$$

- ▶ Greedy approach:
  - ▶ Find best rank-1 approximation:

$$\hat{\mathbf{v}} := \underset{\|\mathbf{x}\|=1}{\arg \max} \mathbf{T}(\mathbf{x}, \mathbf{x}, \mathbf{x}), \quad \hat{\lambda} := \mathbf{T}(\hat{\mathbf{v}}, \hat{\mathbf{v}}, \hat{\mathbf{v}}).$$

- ▶ “Deflate”  $\mathbf{T} := \mathbf{T} - \hat{\lambda} \cdot \hat{\mathbf{v}} \otimes \hat{\mathbf{v}} \otimes \hat{\mathbf{v}}$  and repeat.

## Errors from deflation

(For simplicity, assume  $\lambda_t = 1$  for all  $t$ , so  $\mathbf{T} = \sum_t \mathbf{v}_t^{\otimes 3}$ .)

**First greedy step:**

Rank-1 approx.  $\hat{\mathbf{v}}_1^{\otimes 3}$  to  $\hat{\mathbf{T}}$  satisfies  $\|\hat{\mathbf{v}}_1 - \mathbf{v}_1\| \leq \varepsilon$  (say).

## Errors from deflation

(For simplicity, assume  $\lambda_t = 1$  for all  $t$ , so  $\mathbf{T} = \sum_t \mathbf{v}_t^{\otimes 3}$ .)

**First greedy step:**

Rank-1 approx.  $\hat{\mathbf{v}}_1^{\otimes 3}$  to  $\hat{\mathbf{T}}$  satisfies  $\|\hat{\mathbf{v}}_1 - \mathbf{v}_1\| \leq \varepsilon$  (say).

**Deflation:** To find next  $\mathbf{v}_t$ , use

$$\begin{aligned}\hat{\mathbf{T}} - \hat{\mathbf{v}}_1^{\otimes 3} &= \mathbf{T} + \mathbf{E} - \hat{\mathbf{v}}_1^{\otimes 3} \\ &= \sum_{t=2}^d \mathbf{v}_t^{\otimes 3} + \mathbf{E} + (\mathbf{v}_1^{\otimes 3} - \hat{\mathbf{v}}_1^{\otimes 3}).\end{aligned}$$

## Errors from deflation

(For simplicity, assume  $\lambda_t = 1$  for all  $t$ , so  $\mathbf{T} = \sum_t \mathbf{v}_t^{\otimes 3}$ .)

**First greedy step:**

Rank-1 approx.  $\hat{\mathbf{v}}_1^{\otimes 3}$  to  $\hat{\mathbf{T}}$  satisfies  $\|\hat{\mathbf{v}}_1 - \mathbf{v}_1\| \leq \varepsilon$  (say).

**Deflation:** To find next  $\mathbf{v}_t$ , use

$$\begin{aligned}\hat{\mathbf{T}} - \hat{\mathbf{v}}_1^{\otimes 3} &= \mathbf{T} + \mathbf{E} - \hat{\mathbf{v}}_1^{\otimes 3} \\ &= \sum_{t=2}^d \mathbf{v}_t^{\otimes 3} + \mathbf{E} + (\mathbf{v}_1^{\otimes 3} - \hat{\mathbf{v}}_1^{\otimes 3}).\end{aligned}$$

Now error seems to have **doubled** (i.e., of size  $2\varepsilon$ ) ...

## Effect of deflation errors

For any unit vector  $x$  orthogonal to  $v_1$ :

$$\left\| \frac{1}{3} \nabla_x \left\{ \left( v_1^{\otimes 3} - \hat{v}_1^{\otimes 3} \right) (x, x, x) \right\} \right\| = \left\| \langle v_1, x \rangle^2 v_1 - \langle \hat{v}_1, x \rangle^2 \hat{v}_1 \right\|$$

## Effect of deflation errors

For any unit vector  $x$  orthogonal to  $v_1$ :

$$\begin{aligned}\left\| \frac{1}{3} \nabla_x \left\{ \left( v_1^{\otimes 3} - \hat{v}_1^{\otimes 3} \right) (x, x, x) \right\} \right\| &= \left\| \langle v_1, x \rangle^2 v_1 - \langle \hat{v}_1, x \rangle^2 \hat{v}_1 \right\| \\ &= \langle \hat{v}_1, x \rangle^2\end{aligned}$$

## Effect of deflation errors

For any unit vector  $x$  orthogonal to  $v_1$ :

$$\begin{aligned}\left\| \frac{1}{3} \nabla_x \left\{ \left( v_1^{\otimes 3} - \hat{v}_1^{\otimes 3} \right) (x, x, x) \right\} \right\| &= \left\| \langle v_1, x \rangle^2 v_1 - \langle \hat{v}_1, x \rangle^2 \hat{v}_1 \right\| \\ &= \langle \hat{v}_1, x \rangle^2 \\ &\leq \|v_1 - \hat{v}_1\|^2 \leq \varepsilon^2.\end{aligned}$$

## Effect of deflation errors

For any unit vector  $x$  orthogonal to  $v_1$ :

$$\begin{aligned}\left\| \frac{1}{3} \nabla_x \left\{ \left( v_1^{\otimes 3} - \hat{v}_1^{\otimes 3} \right) (x, x, x) \right\} \right\| &= \left\| \langle v_1, x \rangle^2 v_1 - \langle \hat{v}_1, x \rangle^2 \hat{v}_1 \right\| \\ &= \langle \hat{v}_1, x \rangle^2 \\ &\leq \|v_1 - \hat{v}_1\|^2 \leq \varepsilon^2.\end{aligned}$$

So effect of errors (original and from deflation)  $E + (v_1^{\otimes 3} - \hat{v}_1^{\otimes 3})$  in directions orthogonal to  $v_1$  is  $(1 + o(1))\varepsilon$  rather than  $2\varepsilon$ .

## Effect of deflation errors

For any unit vector  $x$  orthogonal to  $v_1$ :

$$\begin{aligned}\left\| \frac{1}{3} \nabla_x \left\{ \left( v_1^{\otimes 3} - \hat{v}_1^{\otimes 3} \right) (x, x, x) \right\} \right\| &= \left\| \langle v_1, x \rangle^2 v_1 - \langle \hat{v}_1, x \rangle^2 \hat{v}_1 \right\| \\ &= \langle \hat{v}_1, x \rangle^2 \\ &\leq \|v_1 - \hat{v}_1\|^2 \leq \varepsilon^2.\end{aligned}$$

So effect of errors (original and from deflation)  $E + (v_1^{\otimes 3} - \hat{v}_1^{\otimes 3})$  in directions orthogonal to  $v_1$  is  $(1 + o(1))\varepsilon$  rather than  $2\varepsilon$ .

- ▶ Deflation errors have lower-order effect on finding other  $v_t$ .  
(Analogous statement for deflation with matrices does not hold.)

## Summary

- ▶ Using method-of-moments with **low-order moments**, can efficiently **estimate parameters** for many 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**.

# Summary

- ▶ Using method-of-moments with **low-order moments**, can efficiently **estimate parameters** for many 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**.
- ▶ Many issues to resolve!
  - ▶ Handle model misspecification, increase robustness.
  - ▶ General methodology.
  - ▶ Incorporate general prior knowledge and interactive feedback.

# Acknowledgements

**Collaborators:** Anima Anandkumar (Caltech), Rishabh Dudeja (Columbia),  
Dean Foster (Amazon), Rong Ge (Duke), Don Goldfarb (Columbia),  
Sham Kakade (UW), Percy Liang (Stanford), Yi-Kai Liu (NIST),  
Cun Mu (Jet), Matus Telgarsky (UIUC), Tong Zhang (Tencent)

Further reading:

- ▶ [Anandkumar, Ge, H., Kakade, & Telgarsky.](#)

**Tensor decompositions for learning latent variable models.**

*Journal of Machine Learning Research*, 15(Aug):2773–2831, 2014.

<https://goo.gl/F8HudN>



¡Gracias!