



MathSci.ai

Tensor Moments of Gaussian Mixture Models: Theory and Applications

Tammy Kolda
MathSci.ai

Joe Kileel and João M. Pereira
University of Texas, Austin



1/31/2022

Kolda @ Texas A&M TRIPODS Distinguished Lecture

Amazing Coauthors



MathSci.ai



Joe Kileel
Asst Professor
Math Dept
U Texas, Austin



João M. Pereira

Postdoc
Math Dept
U Texas, Austin

(starting summer 2022)
Asst Professor
Instituto Nacional de
Matemática Pura e
Aplicada (IMPA)
Brazil



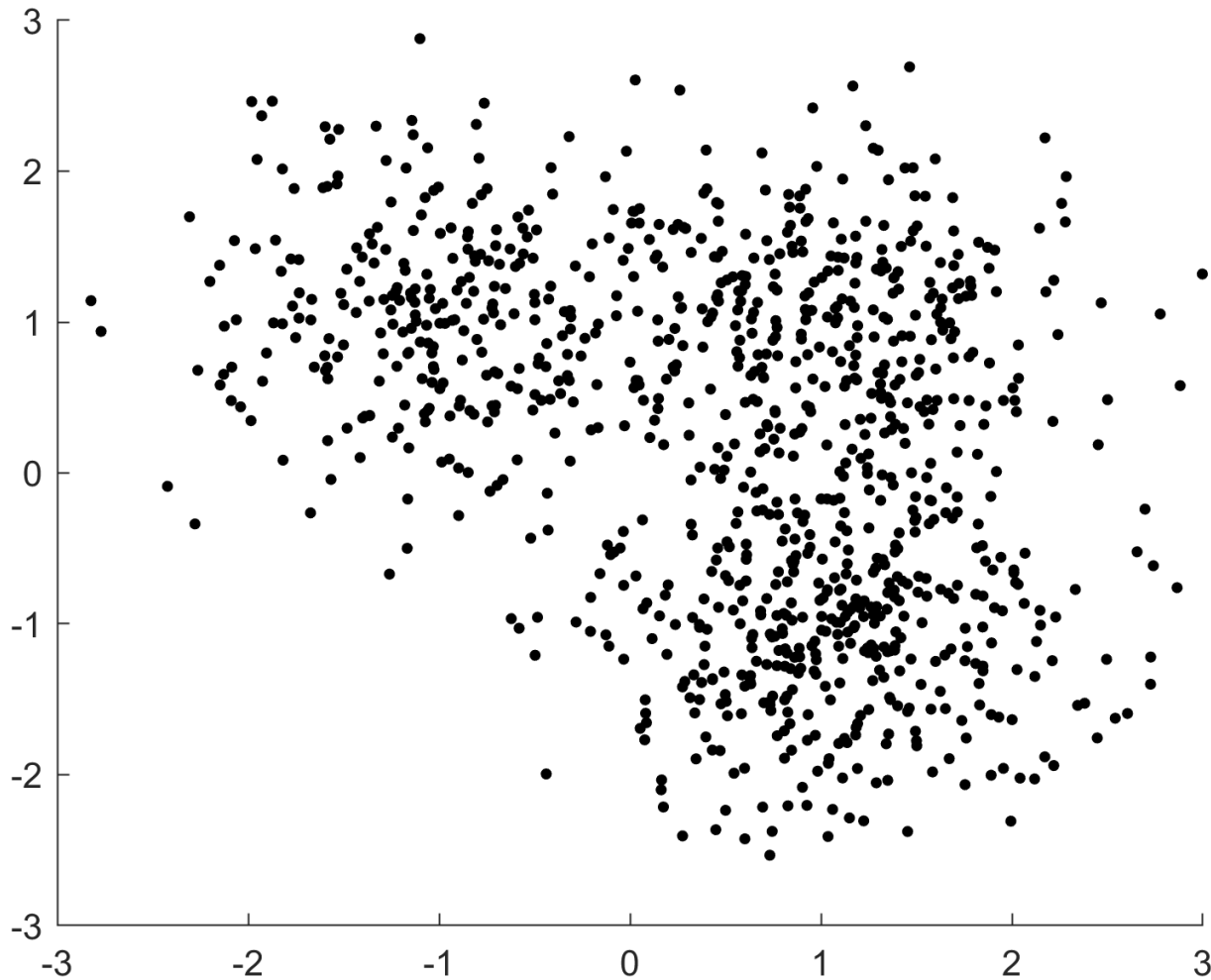
Motivation: Making Sense of Data via Models



Motivation: Making Sense of Data



MathSci.ai

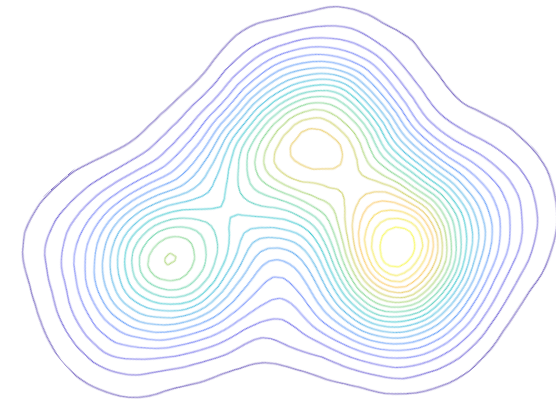
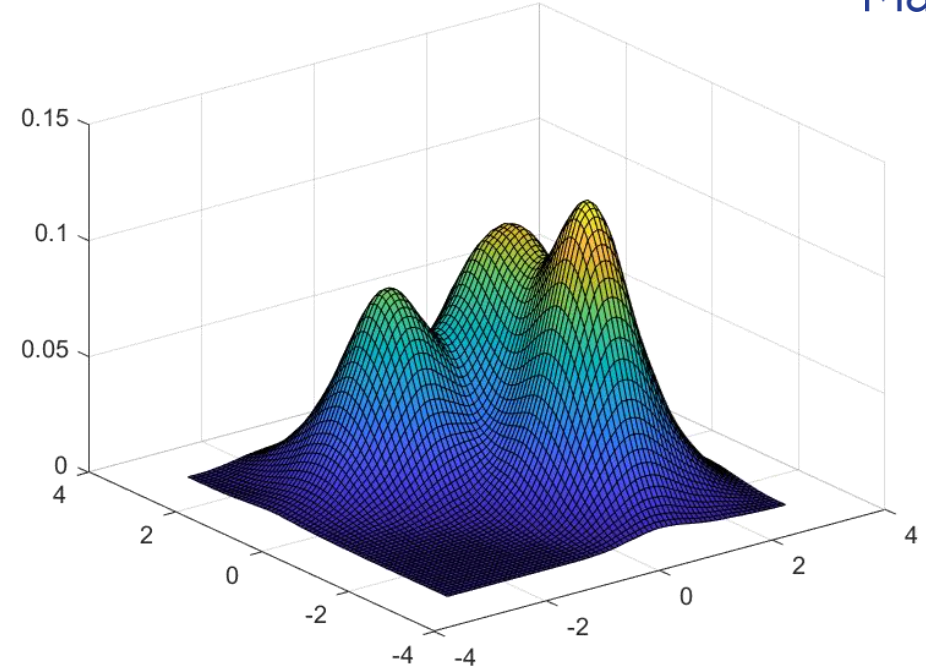
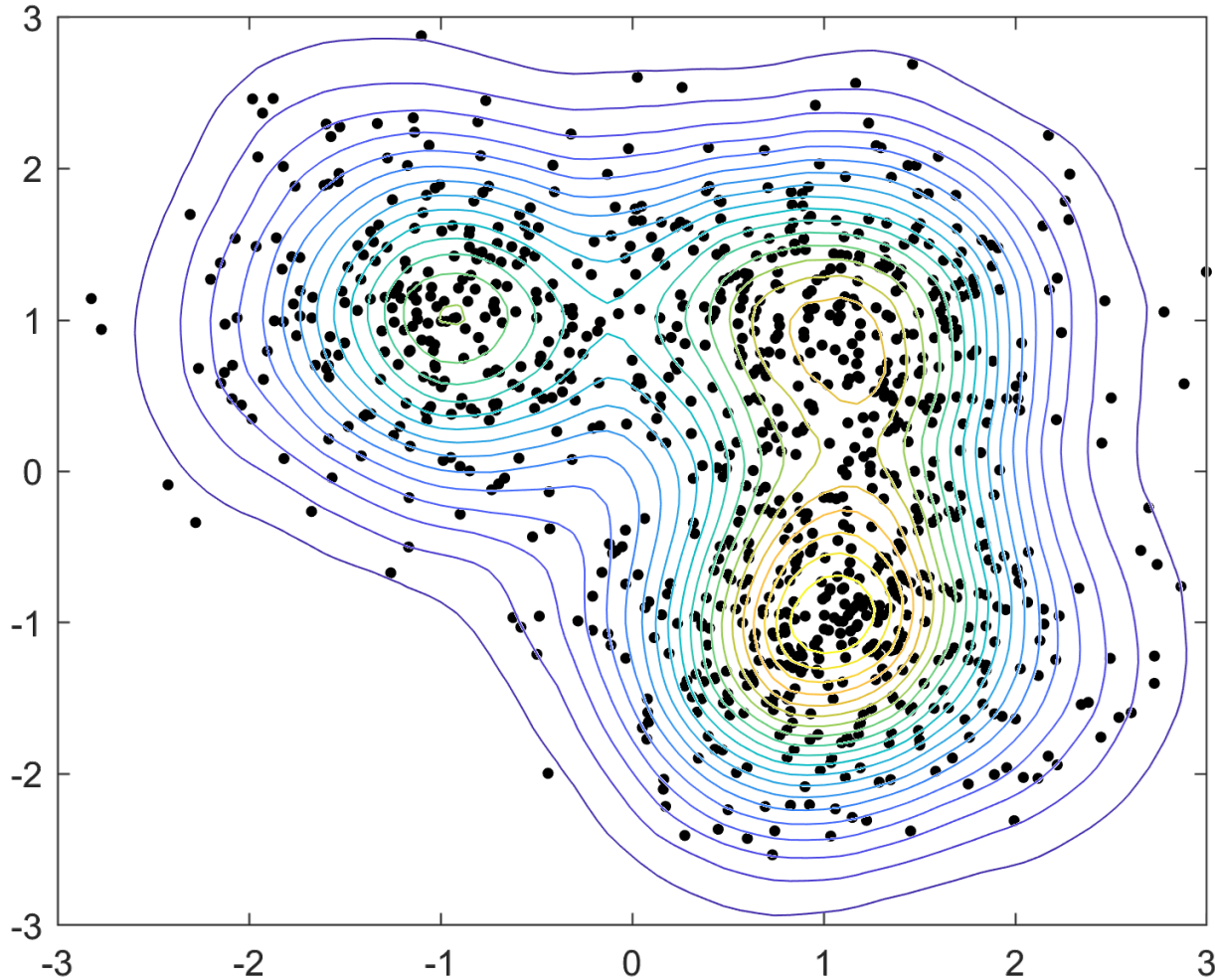


Illustrative Example:
1000 datapoints
 (x_1, x_2)

Technique: Create Model of Probability Density from Datapoints



MathSci.ai





Gaussian Distribution (“Bell curve”)

Probability Distribution Function (pdf)

$$\frac{e^{-\frac{1}{2}(\mathbf{x}-\boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\mathbf{x}-\boldsymbol{\mu})}}{\sqrt{(2\pi)^n |\boldsymbol{\Sigma}|}}$$

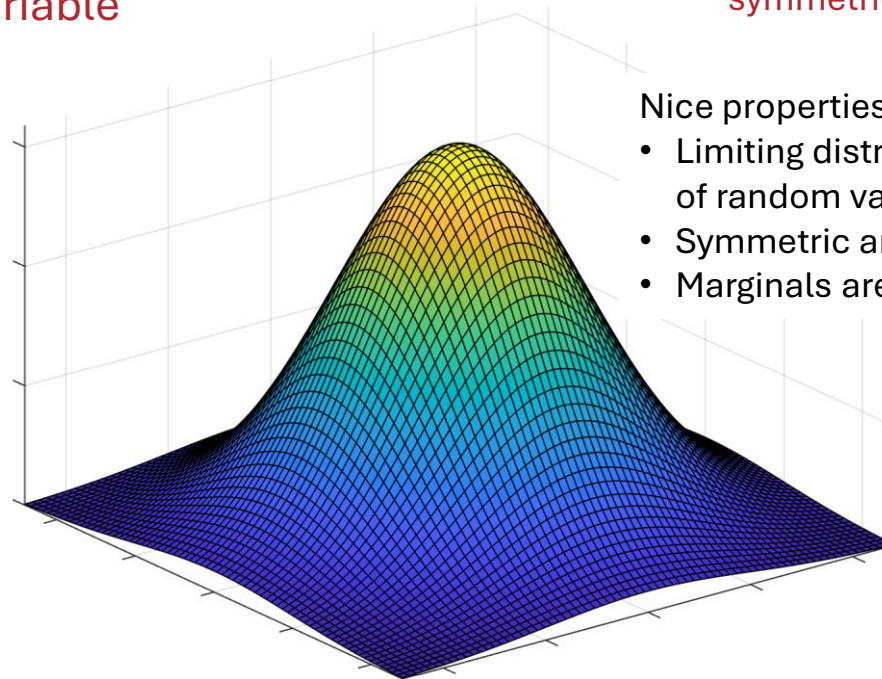
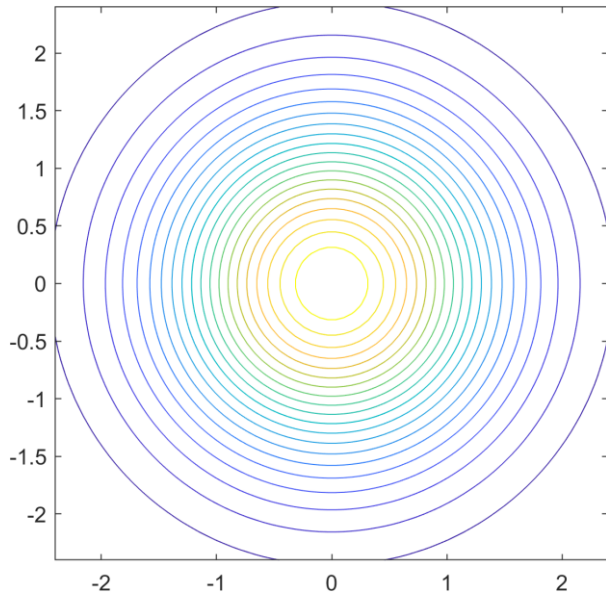
$$X \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$$

X
n-dimensional random variable

$\boldsymbol{\mu}$
mean or center
 $\boldsymbol{\mu} \in \mathbb{R}^n$

$\boldsymbol{\Sigma}$
covariance matrix
 $\boldsymbol{\Sigma} \in \mathbb{R}^{n \times n}$
symmetric positive definite

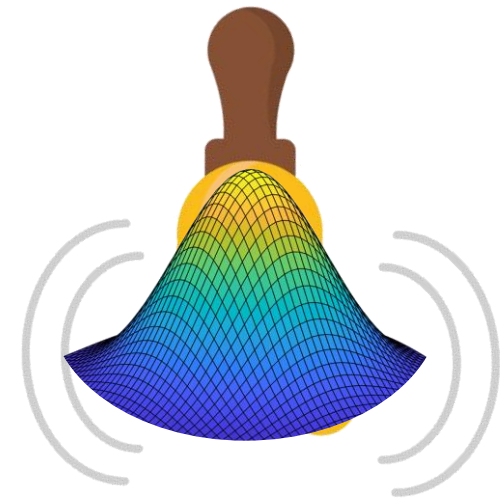
aka “Normal” distribution



Nice properties...

- Limiting distribution of a sum of random variables
- Symmetric around the mean
- Marginals are normal

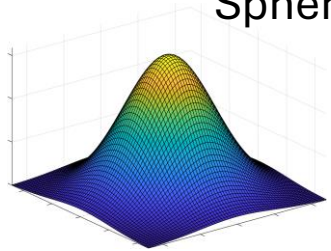
Bell Curve



Covariance for Gaussian Distribution

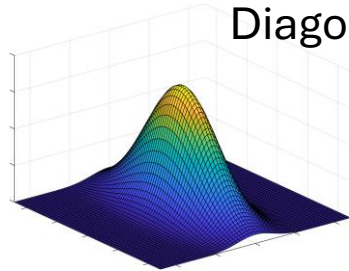


Spherical or Isotropic



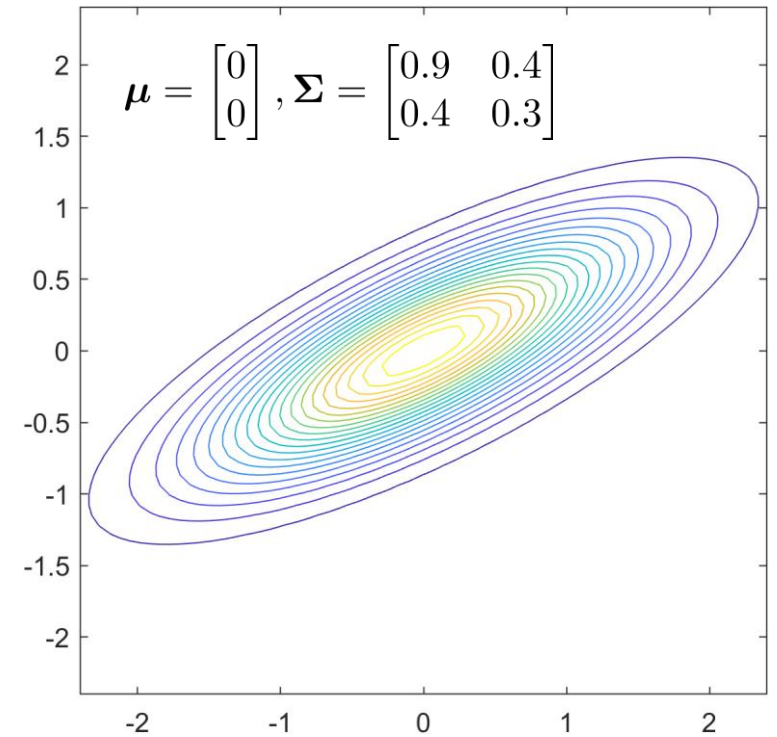
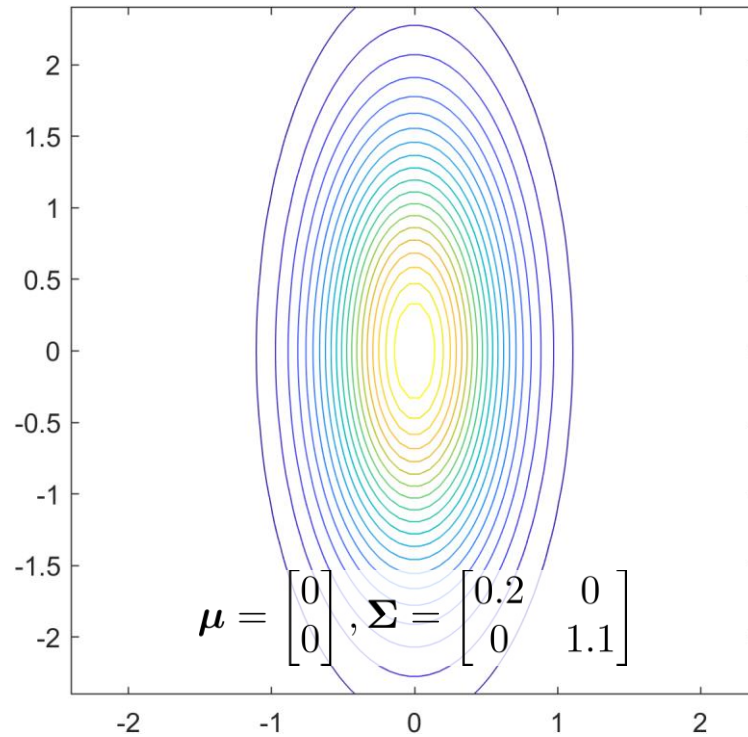
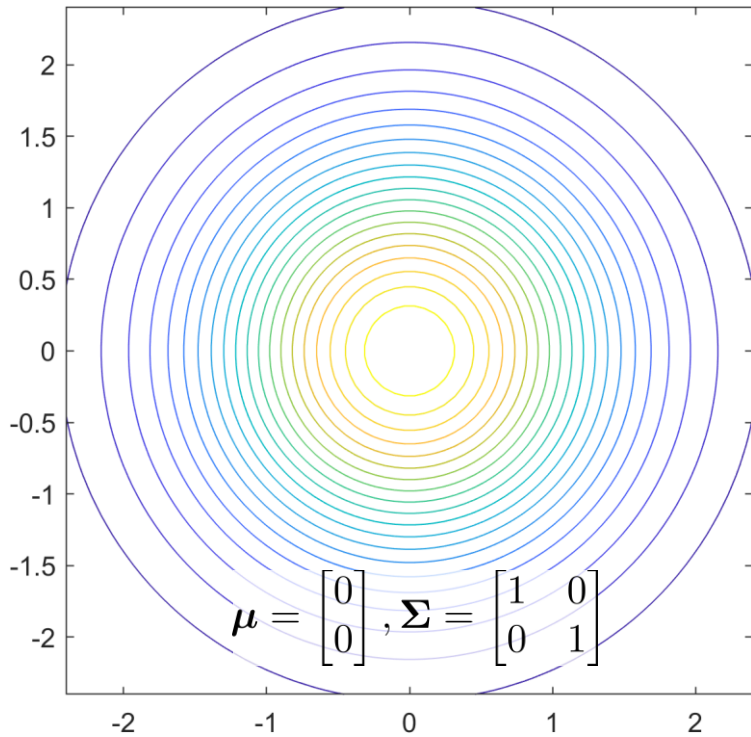
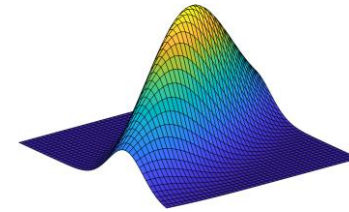
$$\Sigma = \sigma^2 \mathbf{I}$$

Diagonal or Axis-Aligned

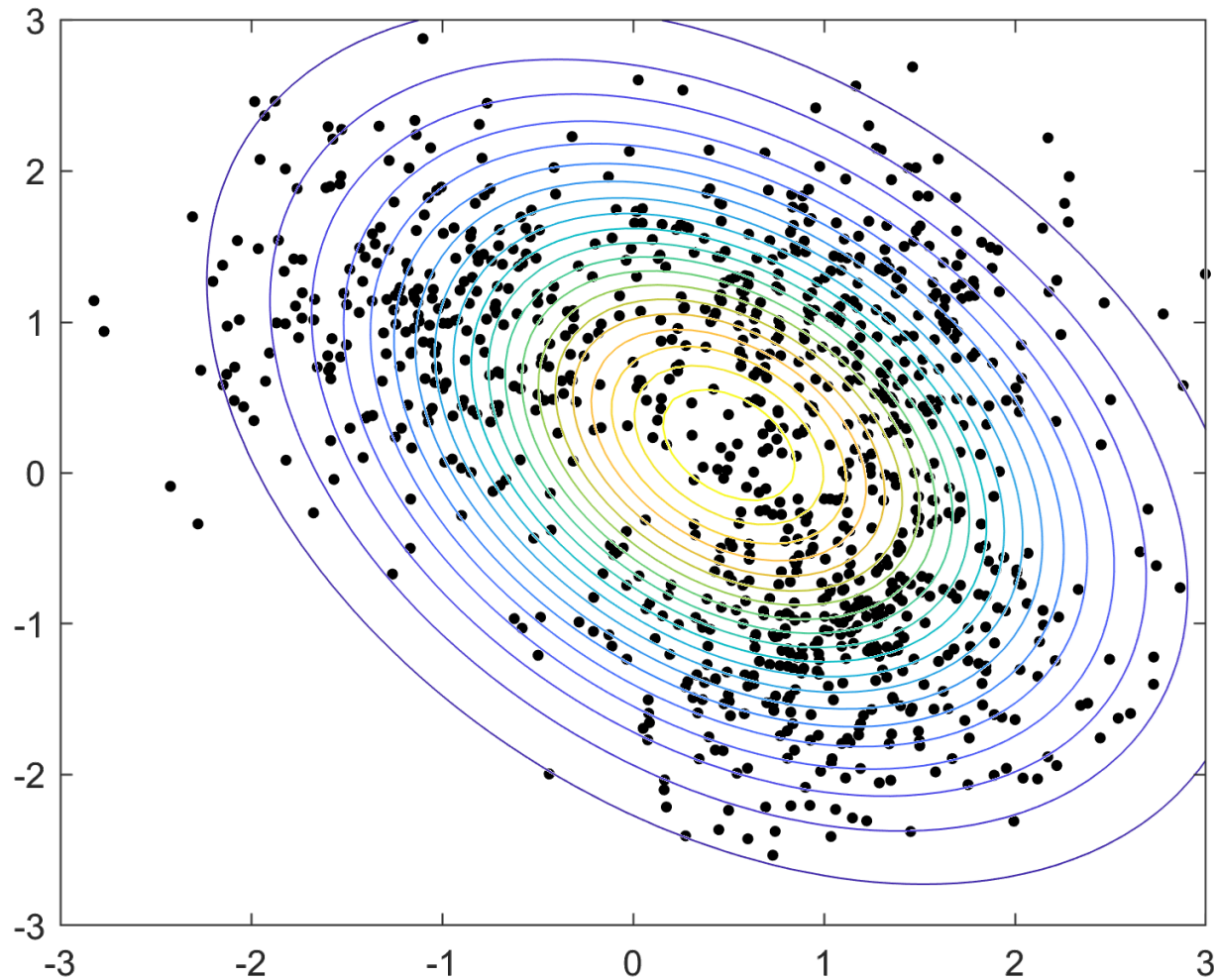


$$\Sigma = \begin{bmatrix} \sigma_1^2 & 0 \\ 0 & \sigma_2^2 \end{bmatrix}$$

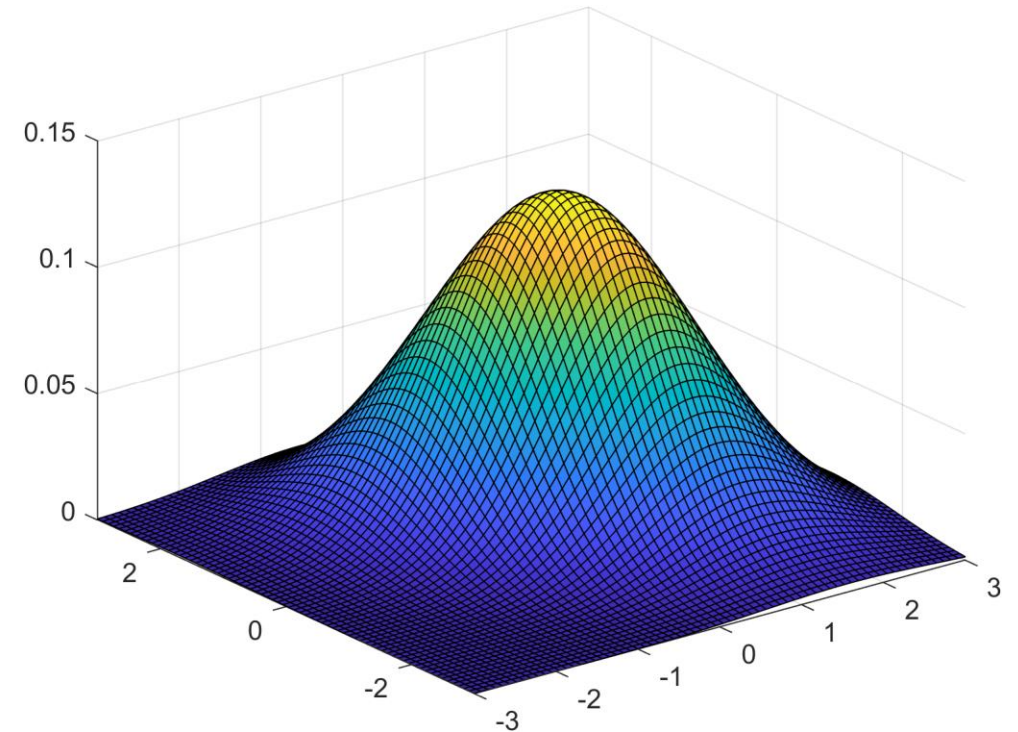
General



Gaussian Model of Data



Parameters: μ, Σ



Gaussian Mixture Model (GMM)

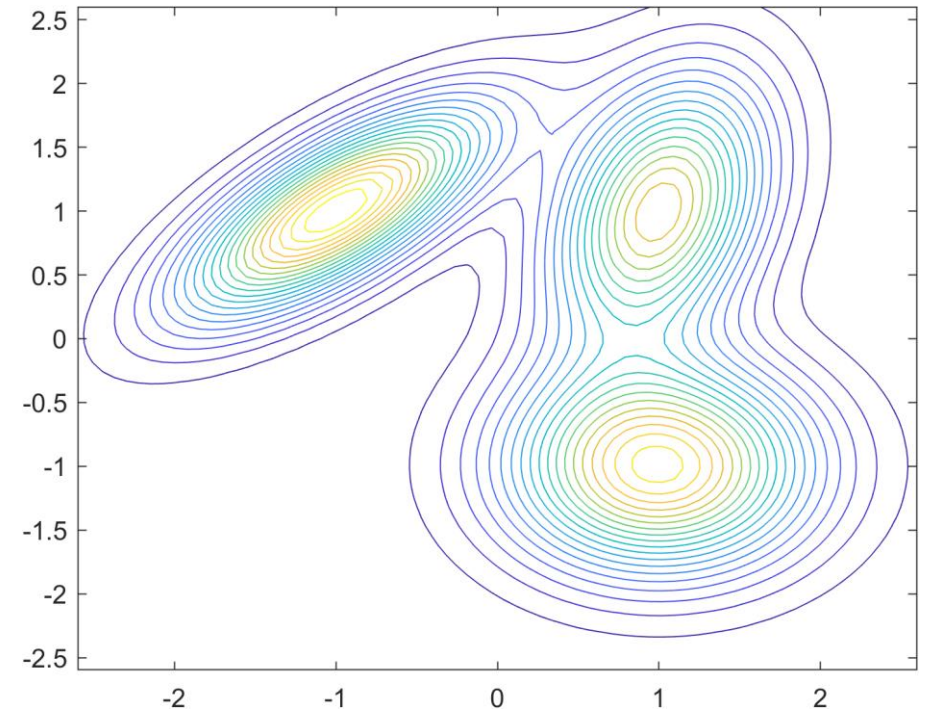
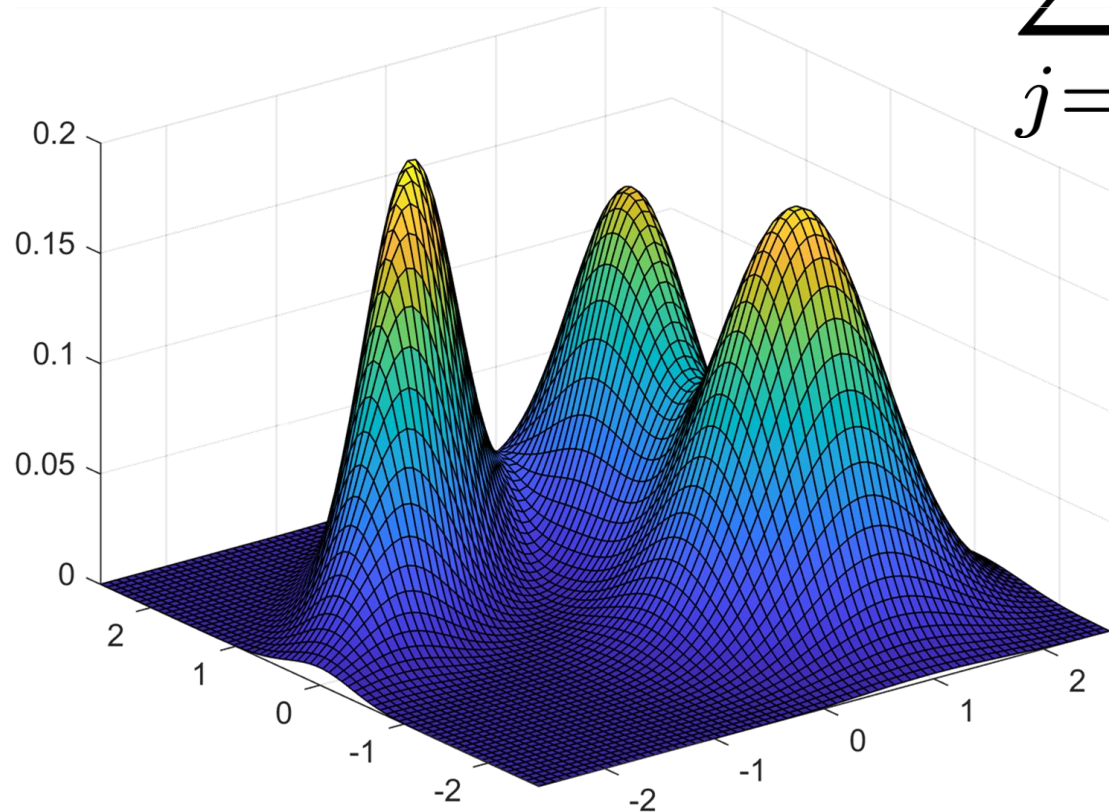


Number of Components $\rightarrow m$

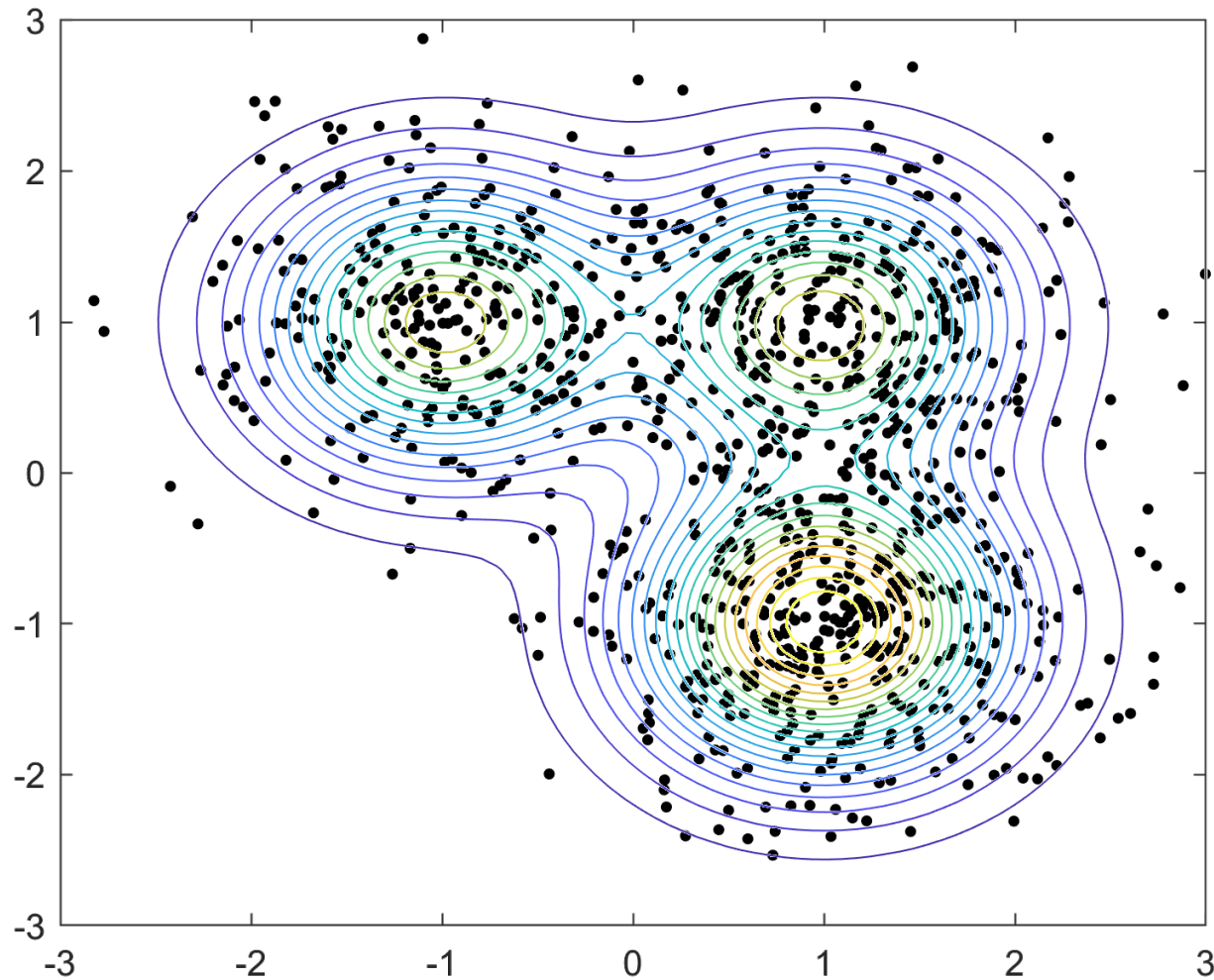
Component j

$$X \sim \sum_{j=1}^m \lambda_j \mathcal{N}(\mu_j, \Sigma_j)$$

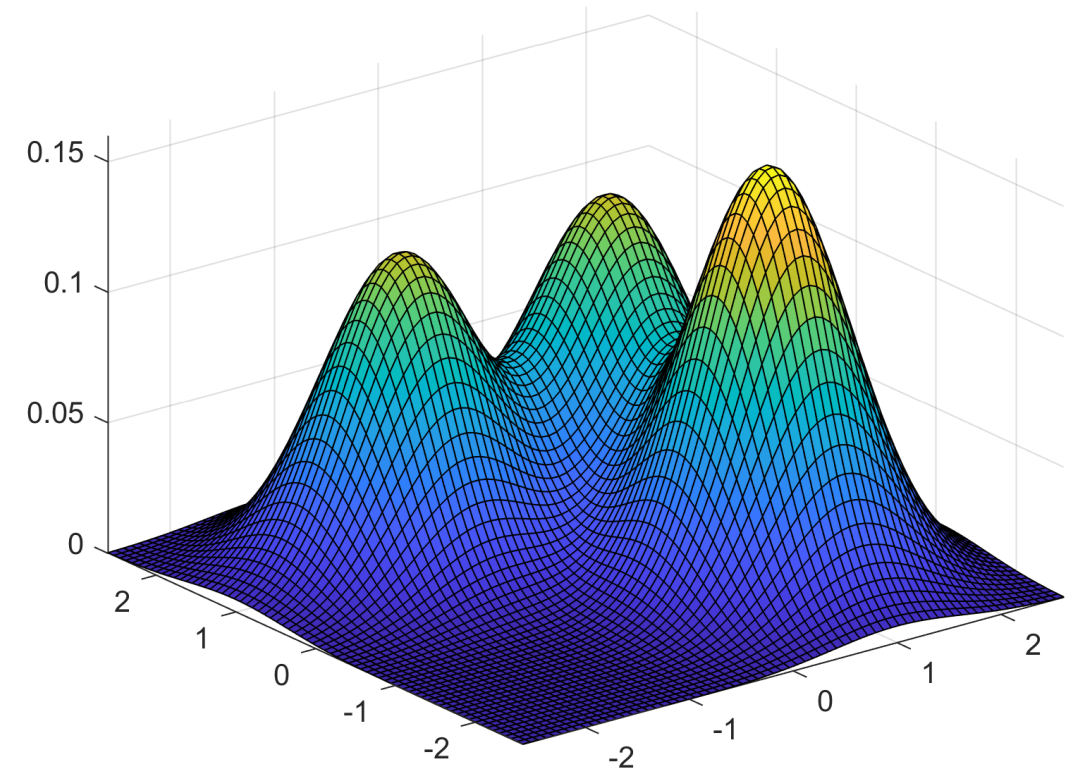
Proportion in j th component



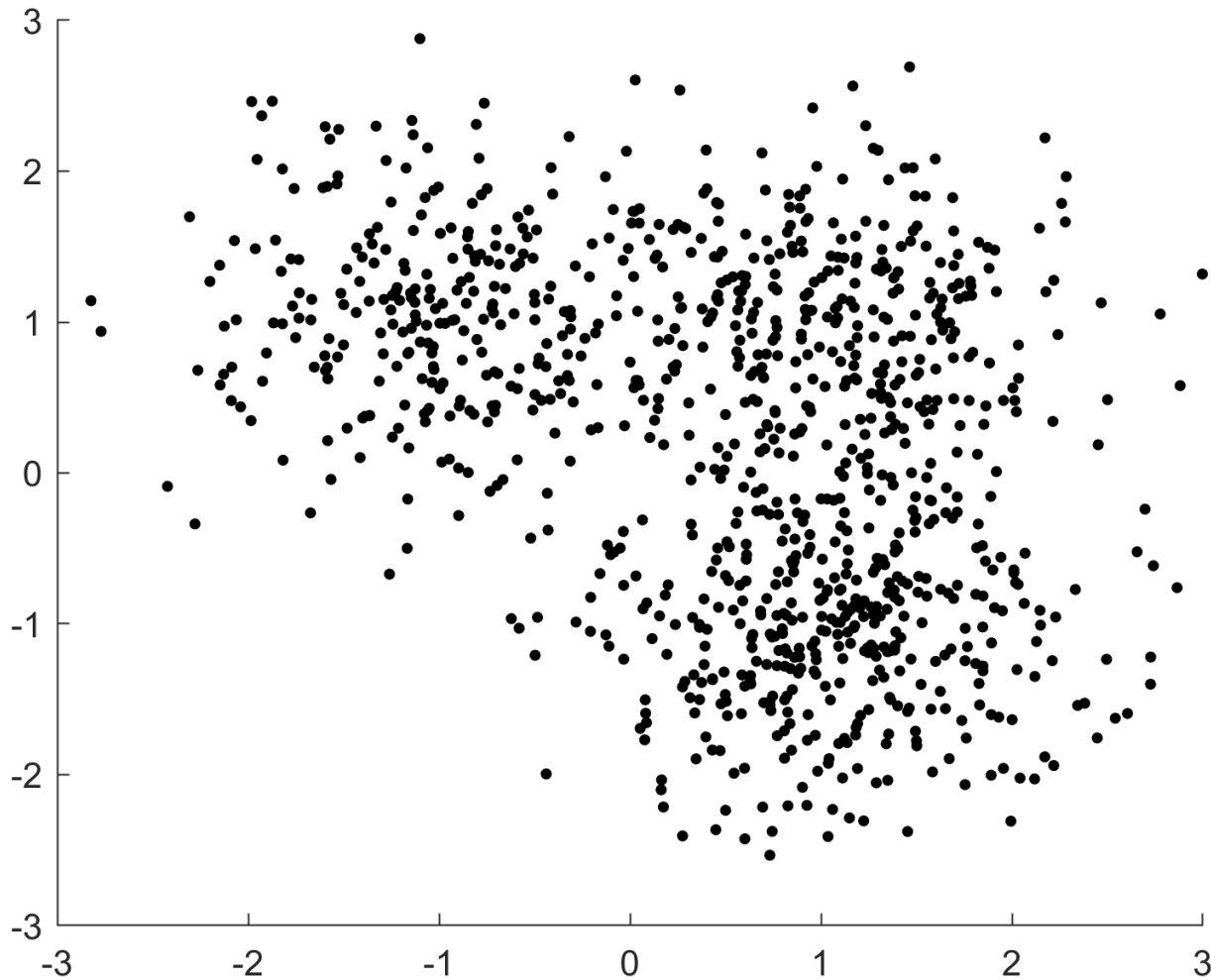
GMM with Three Components



Parameters: $\{\lambda_j, \mu_j, \Sigma_j\}_{j=1:m}$



Motivation: Making Sense of Data



Illustrative Example:

1000 datapoints

(x_1, x_2)

Real-World Scenario:

10,000+ datapoints

$(x_1, x_2, \dots, x_{500})$

Many more
datapoints

Much higher
dimension

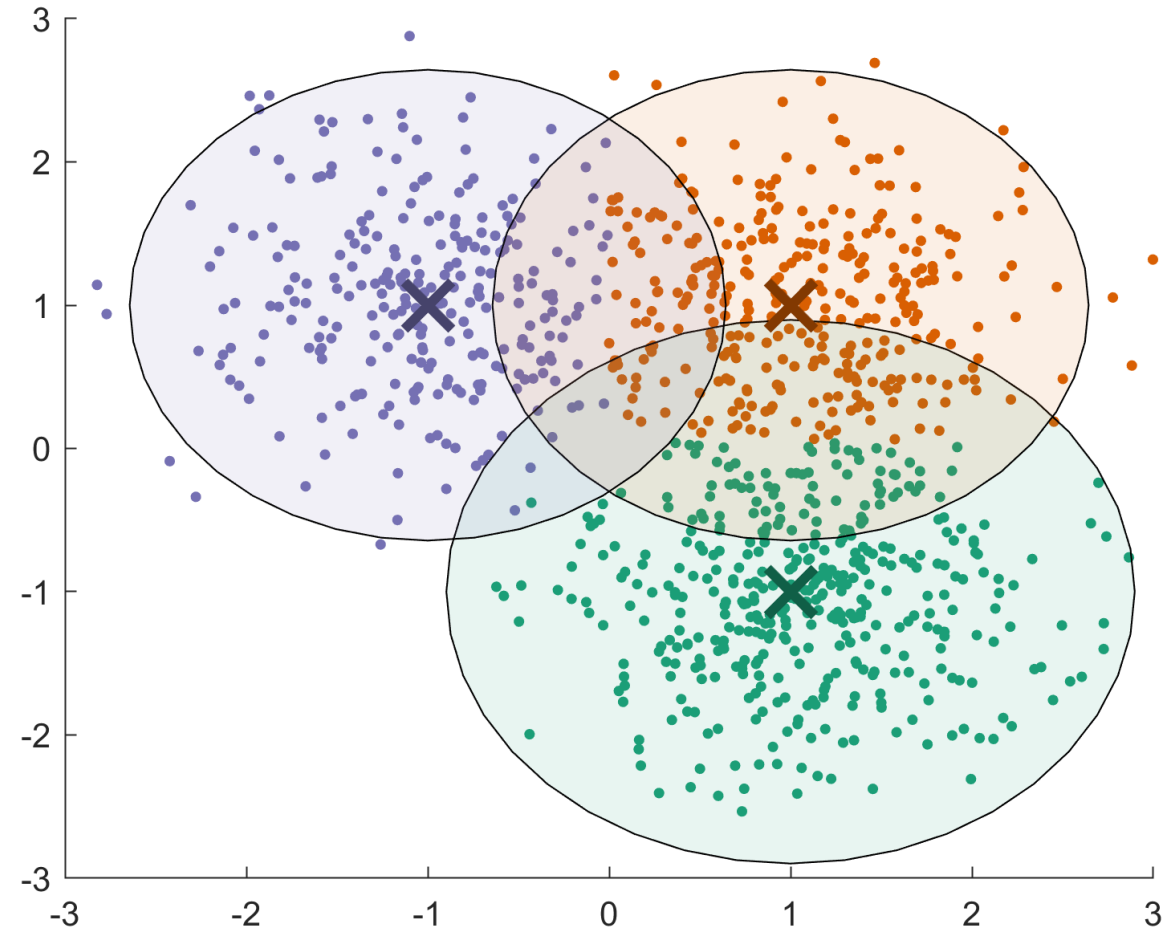
Applications: Gaussian Mixtures



Density Estimation

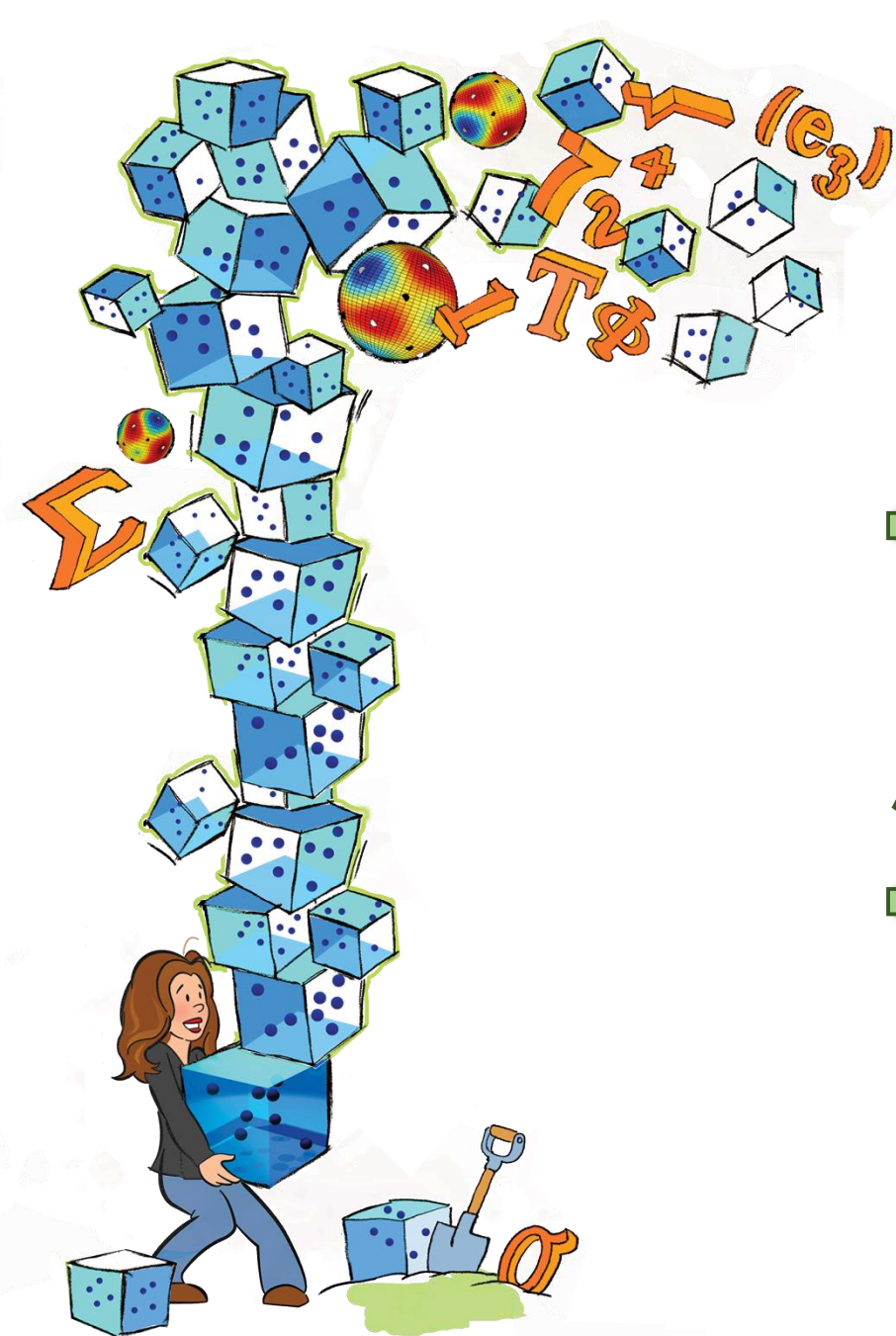
Clustering

Anomaly Detection





Tensors and Symmetric Tensors



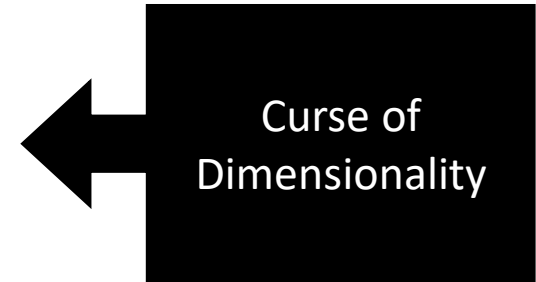
Tensors are Multi-dimensional Arrays



$d = \text{order}$ of the tensor (the number of ways or modes)

For this talk, all modes have the same size: n

$n^d =$ number of entries for d -way tensor of dimension n



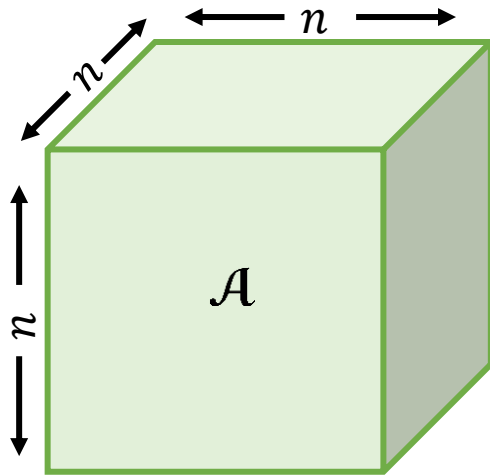
$a_{ijk} = (i, j, k)$ entry of 3-way tensor \mathcal{A}

A tensor is **symmetric** if its entries are invariant under permutation, i.e.,

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

Curse of notation...

$(i_1, i_2, \dots, i_d) = \text{index}$ into tensor, $i_k \in \{1, \dots, n\}$ for $k = 1, 2, \dots, d$

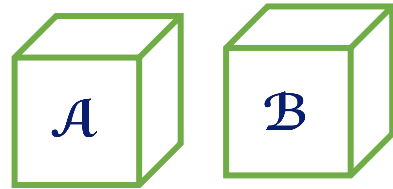


3-way tensor

Tensor Norm & Inner Product

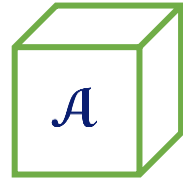


Inner Product

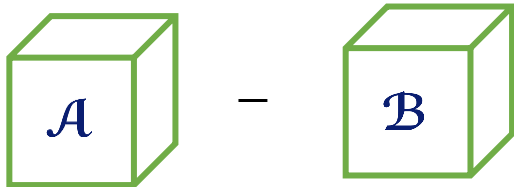


$$\langle \mathcal{A}, \mathcal{B} \rangle = \sum_{i=1}^n \sum_{j=1}^n \sum_{k=1}^n a_{ijk} b_{ijk}$$

Norm



$$\|\mathcal{A}\|^2 = \langle \mathcal{A}, \mathcal{A} \rangle = \sum_{i=1}^n \sum_{j=1}^n \sum_{k=1}^n a_{ijk}^2$$



$$\|\mathcal{A} - \mathcal{B}\|^2 = \|\mathcal{A}\|^2 + \|\mathcal{B}\|^2 - 2\langle \mathcal{A}, \mathcal{B} \rangle$$

Symmetric Tensor Outer Product



$$\mathbf{b} \in \mathbb{R}^n$$

$$\mathbf{A} = \mathbf{b}^{\otimes 2} = \mathbf{b} \otimes \mathbf{b} \in \mathbb{R}^{n \times n}$$



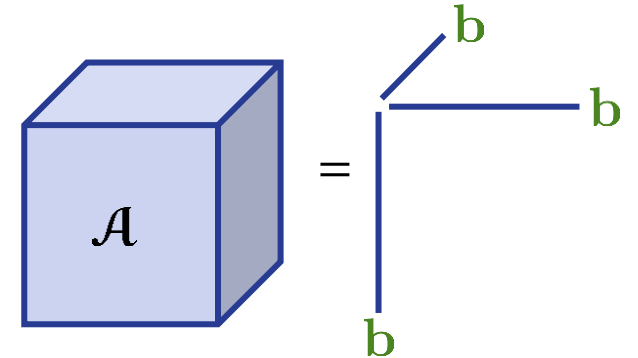
$$a_{ij} = b_i b_j$$

$$\mathcal{A} = \mathbf{b}^{\otimes 3} = \mathbf{b} \otimes \mathbf{b} \otimes \mathbf{b} \in \mathbb{R}^{n \times n \times n}$$



$$a_{ijk} = b_i b_j b_k$$

Visualization of 3-way
Outer Product



$$\mathcal{A} = \mathbf{b}^{\otimes 4} = \mathbf{b} \otimes \mathbf{b} \otimes \mathbf{b} \otimes \mathbf{b} \in \mathbb{R}^{n \times n \times n \times n}$$



$$a_{ijkl} = b_i b_j b_k b_l$$

$$\mathcal{A} = \mathbf{b}^{\otimes d} = \underbrace{\mathbf{b} \otimes \dots \otimes \mathbf{b}}_{d \text{ times}} \in \mathbb{R}^{n \times \dots \times n}$$

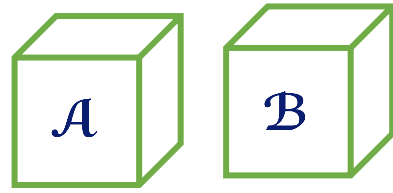


$$a_{i_1 \dots i_d} = b_{i_1} \dots b_{i_d}$$

Outer Products, Inner Products, and Norms

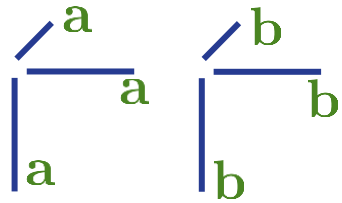


Inner Product



$$\langle \mathcal{A}, \mathcal{B} \rangle = \sum_{i=1}^n \sum_{j=1}^n \sum_{k=1}^n a_{ijk} b_{ijk}$$

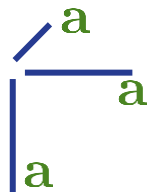
Inner Product of Outer Products



$$\begin{aligned} \langle \mathbf{a}^{\otimes 3}, \mathbf{b}^{\otimes 3} \rangle &= \sum_{i=1}^n \sum_{j=1}^n \sum_{k=1}^n (a_i a_j a_k) (b_i b_j b_k) \\ &= \left(\sum_{i=1}^n a_i b_i \right) \left(\sum_{j=1}^n a_j b_j \right) \left(\sum_{k=1}^n a_k b_k \right) = \langle \mathbf{a}, \mathbf{b} \rangle^3 \end{aligned}$$

$$\langle \mathbf{a}^{\otimes d}, \mathbf{b}^{\otimes d} \rangle = \langle \mathbf{a}, \mathbf{b} \rangle^d$$

Norm

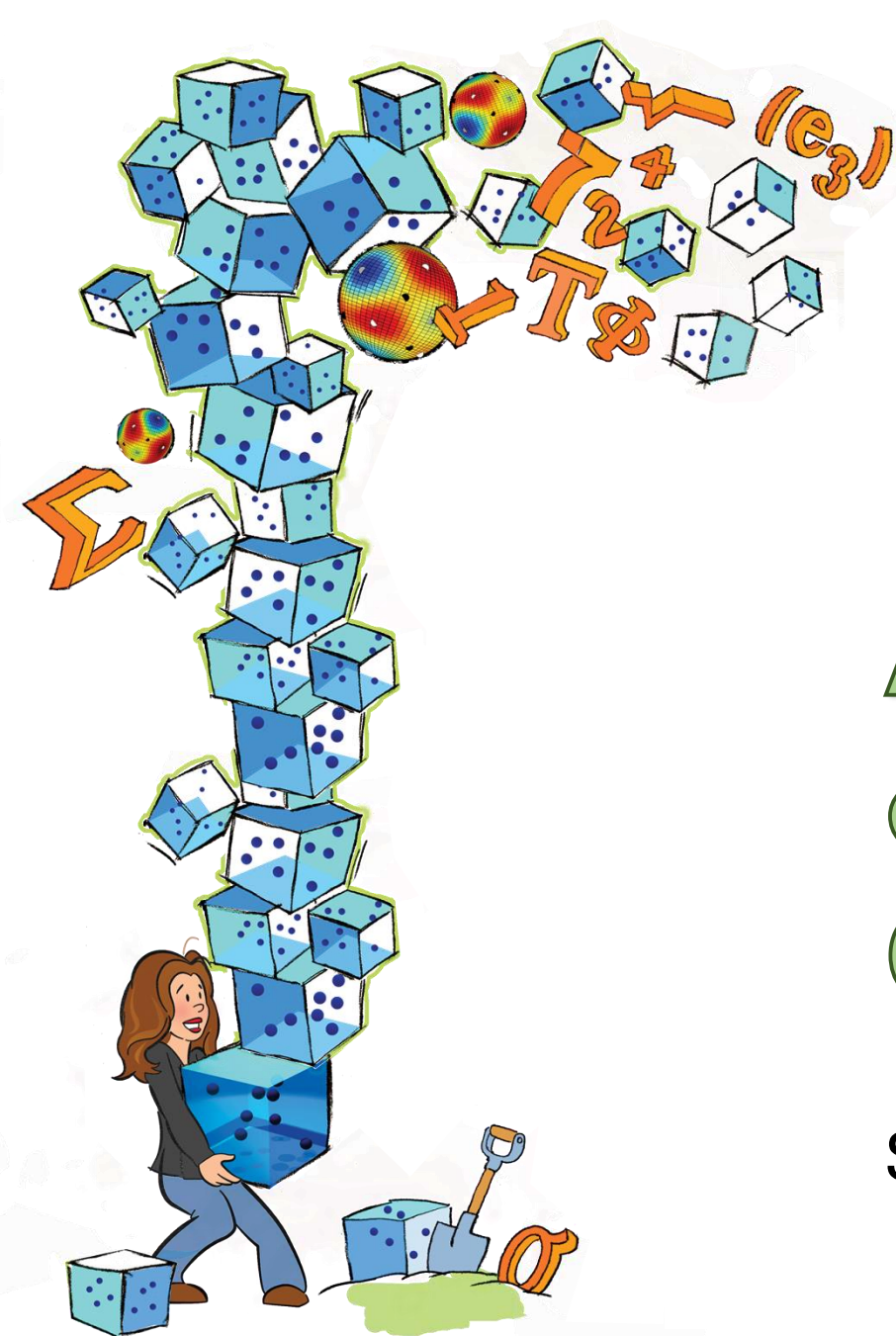


$$\|\mathbf{a}^{\otimes d}\|^2 = \langle \mathbf{a}^{\otimes d}, \mathbf{a}^{\otimes d} \rangle = \langle \mathbf{a}, \mathbf{a} \rangle^d = \|\mathbf{a}\|^{2d}$$



Approximate Method of Moments for Gaussian Mixtures

See Sherman & K (2020)



Moments of a Multivariate Random Variable



Higher moments capture interactions between the variables

Random variable: $X \in \mathbb{R}^n$

First moment: $\mathcal{M}^{(1)} = \mathbb{E}(X)$ $\mathcal{M}^{(1)}(i) = \mathbb{E}(X_i)$

Second moment: $\mathcal{M}^{(2)} = \mathbb{E}(X^{\otimes 2})$ $\mathcal{M}^{(2)}(i, j) = \mathbb{E}(X_i X_j)$

Third moment: $\mathcal{M}^{(3)} = \mathbb{E}(X^{\otimes 3})$ $\mathcal{M}^{(3)}(i, j, k) = \mathbb{E}(X_i X_j X_k)$

d th moment: $\mathcal{M}^{(d)} = \mathbb{E}(X^{\otimes d})$

d th moment is a symmetric tensor of order d

Moments Define a Distribution



Higher moments capture interactions between the variables

Random variable: $X \in \mathbb{R}^n$

First moment: $\mathcal{M}^{(1)} = \mathbb{E}(X)$

Second moment: $\mathcal{M}^{(2)} = \mathbb{E}(X^{\otimes 2})$

Third moment: $\mathcal{M}^{(3)} = \mathbb{E}(X^{\otimes 3})$

d th moment: $\mathcal{M}^{(d)} = \mathbb{E}(X^{\otimes d})$

“Method of Moments”
matches empirical and
model moments to
estimate the parameters
of a distribution.

We focus primarily on matching just the d th moment

Gaussian Mixture Model: Small Spherical Covariance



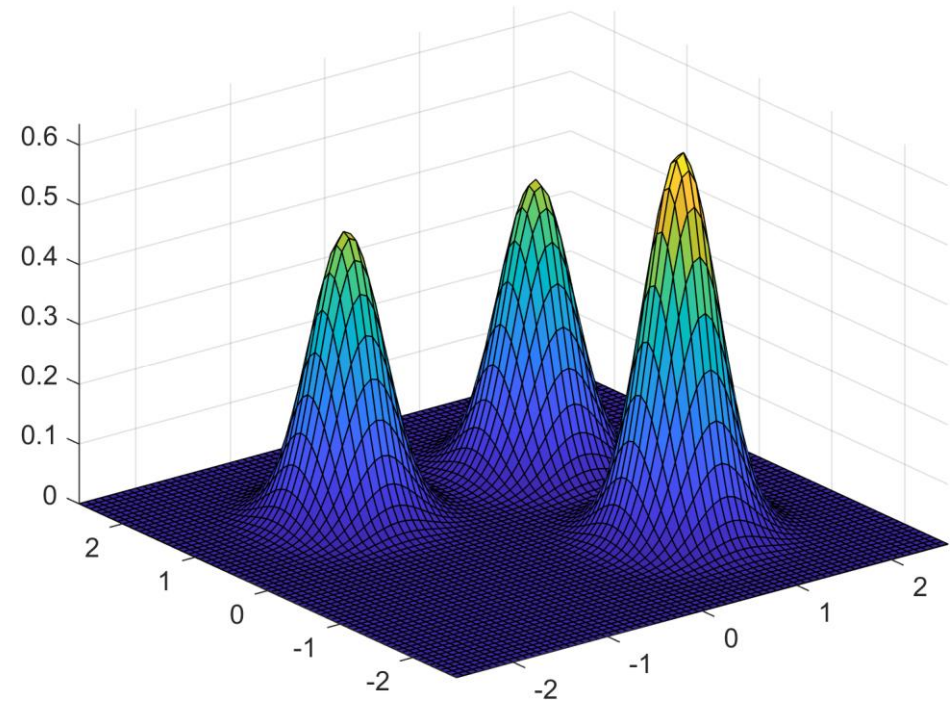
$$X \sim \sum_{j=1}^m \lambda_j \mathcal{N}(\boldsymbol{\mu}_j, \sigma^2 \mathbf{I})$$

Moment of the Model:

$$\mathcal{M}^{(d)} = \sum_{j=1}^m \lambda_j \boldsymbol{\mu}_j^{\otimes d} + \mathcal{O}(\sigma^2)$$

Approximate Moment of the Model:

$$\widetilde{\mathcal{M}}^{(d)} = \sum_{j=1}^m \lambda_j \boldsymbol{\mu}_j^{\otimes d}$$



Fitting the d th Moment

Given r observations:

$$\{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_p\} \subset \mathbb{R}^n$$

Find m weights and mean vectors:

$$\{\lambda_1, \lambda_2, \dots, \lambda_m\} \subset \mathbb{R}$$

$$\{\boldsymbol{\mu}_1, \boldsymbol{\mu}_2, \dots, \boldsymbol{\mu}_m\} \subset \mathbb{R}^n$$

*ignoring
covariance
 $\sigma^2 \mathbf{I}$*

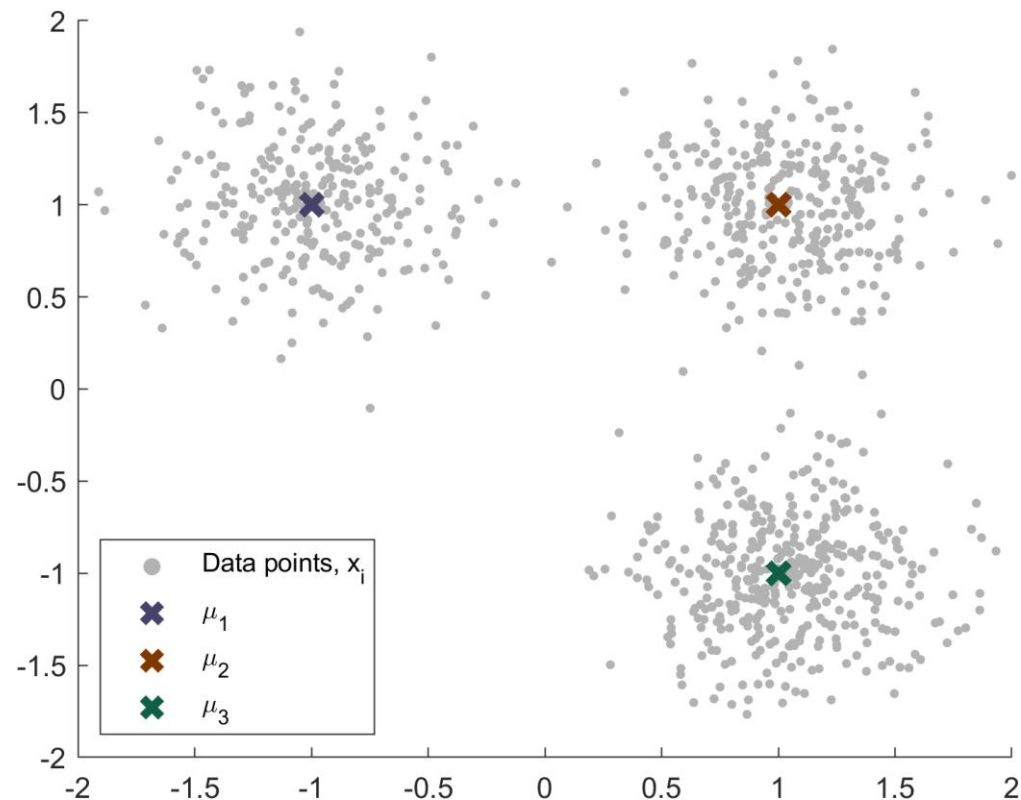
Such that...

$$\sum_{j=1}^m \lambda_j \boldsymbol{\mu}_j^{\otimes d} \approx \frac{1}{p} \sum_{i=1}^p \mathbf{x}_i^{\otimes d}$$

Empirical
or Sample
Moment

Approximate
Model
Moment

Example: $r = 1000, m = 3$



Optimization Formulation: Symmetric Tensor Decomposition

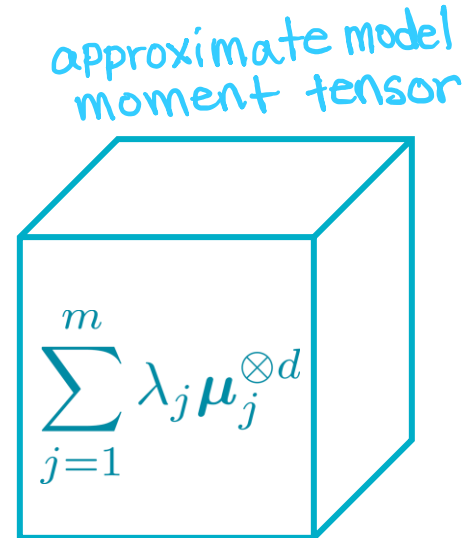
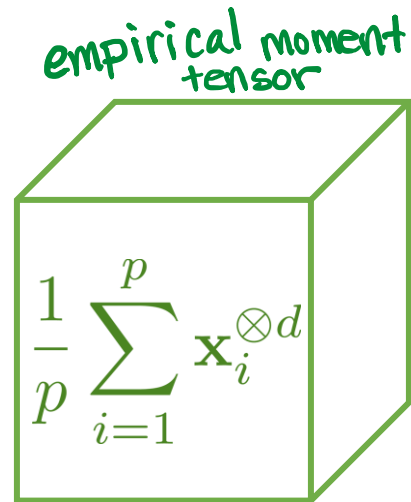


Optimization Parameters

$$\theta = \{\lambda_j, \mu_j\}_{j=1}^m$$

$$\min_{\theta} f(\theta) \equiv \left\| \frac{1}{p} \sum_{i=1}^p \mathbf{x}_i^{\otimes d} - \sum_{j=1}^m \lambda_j \mu_j^{\otimes d} \right\|^2$$

Problem: Forming and storing the empirical moment tensors costs $\mathcal{O}(pn^d)$ operations and $\mathcal{O}(n^d)$ storage



Problem: Forming and storing the approximate model moment tensor costs $\mathcal{O}(mn^d)$ operations and $\mathcal{O}(n^d)$ storage

But there is only $\mathcal{O}(pn)$ data and $\mathcal{O}(mn)$ parameters, so there is room for efficiency

Optimization Formulation as Inner Products



$$\min_{\theta} f(\theta) \equiv \left\| \frac{1}{p} \sum_{i=1}^p \mathbf{x}_i^{\otimes d} - \sum_{j=1}^m \lambda_j \boldsymbol{\mu}_j^{\otimes d} \right\|^2 \quad \theta = \{ \lambda_j, \boldsymbol{\mu}_j \}_{j=1}^m$$

$$f(\theta) = \left\| \frac{1}{p} \sum_{i=1}^p \mathbf{x}_i^{\otimes d} \right\|^2 + \left\| \sum_{j=1}^m \lambda_j \boldsymbol{\mu}_j^{\otimes d} \right\|^2 - 2 \left\langle \frac{1}{p} \sum_{i=1}^p \mathbf{x}_i^{\otimes d}, \sum_{j=1}^m \lambda_j \boldsymbol{\mu}_j^{\otimes d} \right\rangle$$

constant

$$f(\theta) = C + \sum_{i=1}^m \sum_{j=1}^m \lambda_i \lambda_j \langle \boldsymbol{\mu}_i^{\otimes d}, \boldsymbol{\mu}_j^{\otimes d} \rangle - \frac{2}{p} \sum_{i=1}^p \sum_{j=1}^m \lambda_j \langle \mathbf{x}_i^{\otimes d}, \boldsymbol{\mu}_j^{\otimes d} \rangle$$

$$f(\theta) = C + \sum_{i=1}^m \sum_{j=1}^m \lambda_i \lambda_j \langle \boldsymbol{\mu}_i, \boldsymbol{\mu}_j \rangle^d - \frac{2}{p} \sum_{i=1}^p \sum_{j=1}^m \lambda_j \langle \mathbf{x}_i, \boldsymbol{\mu}_j \rangle^d$$

Dot products!



Casting as Optimization Problem

$$\frac{1}{p} \sum_{i=1}^p \mathbf{x}_i^{\otimes d}$$

$$\sum_{j=1}^m \lambda_j \boldsymbol{\mu}_j^{\otimes d}$$

$$f(\theta) = C + \sum_{i=1}^m \sum_{j=1}^m \lambda_i \lambda_j \langle \boldsymbol{\mu}_i, \boldsymbol{\mu}_j \rangle^d - \frac{2}{p} \sum_{i=1}^p \sum_{j=1}^m \lambda_j \langle \mathbf{x}_i, \boldsymbol{\mu}_j \rangle^d$$

$$\theta = \{ \lambda_j, \boldsymbol{\mu}_j \}_{j=1}^m$$

- Never need to form empirical or approximate model moments explicitly, overcoming curse of dimensionality
- Function can be calculated with only dot products, total work $\mathcal{O}(m^2n + pmn)$ and $\mathcal{O}(mn + pn)$ storage, **versus $\mathcal{O}(mn^d + pn^d)$ as originally formulated**
- Gradients equally efficient to calculate, via chain rule
- *Issue:* Inherent scaling problem (will come back to this later)
- Easy stochastic function and gradient if number of samples (p) is large...

$$\tilde{f}(\theta) = C + \sum_{i=1}^m \sum_{j=1}^m \lambda_i \lambda_j \langle \boldsymbol{\mu}_i, \boldsymbol{\mu}_j \rangle^d - \frac{2}{|\Omega|} \sum_{i \in \Omega} \sum_{j=1}^m \lambda_j \langle \mathbf{x}_i, \boldsymbol{\mu}_j \rangle^d$$



But the Approximation is Biased

$$X \sim \sum_{j=1}^m \lambda_j \mathcal{N}(\mu_j, \sigma^2 \mathbf{I})$$

$$\mathbb{E}(X^{\otimes d}) = \sum_{j=1}^m \lambda_j \mu_j^{\otimes d} + \mathcal{O}(\sigma^2)$$

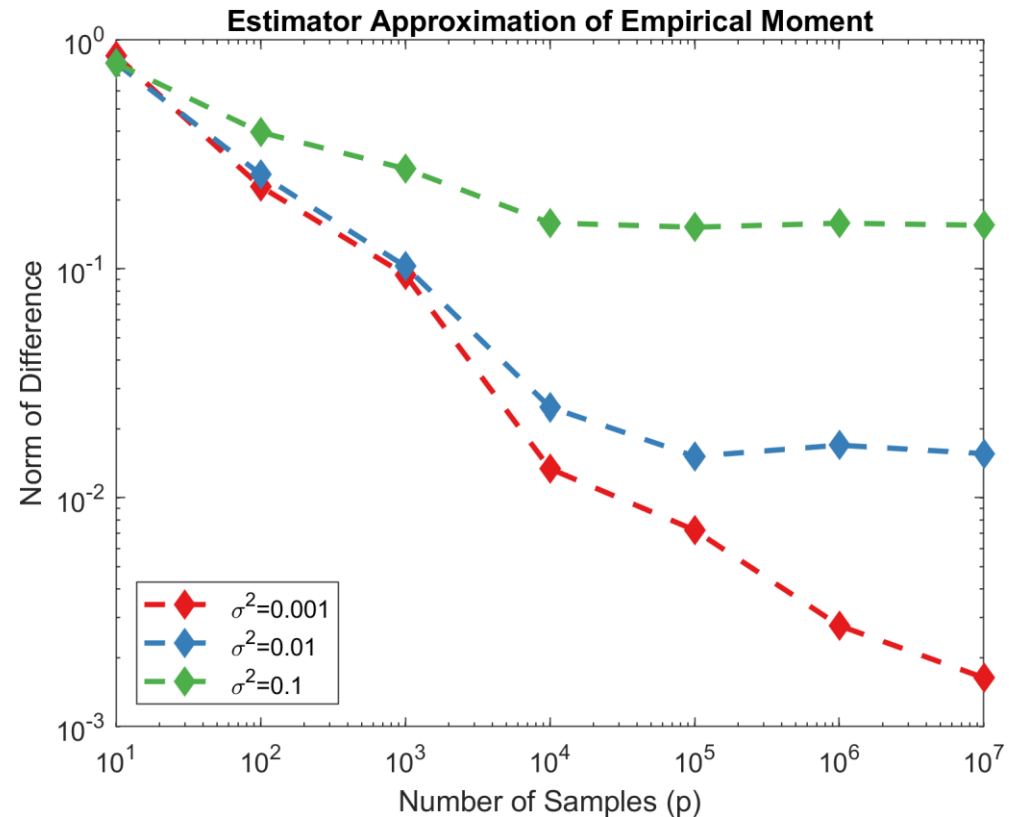
approximate
model
moment
tensor

$$\tilde{\mathcal{M}}^{(d)} = \sum_{j=1}^m \lambda_j \mu_j^{\otimes d}$$

empirical
moment
tensor

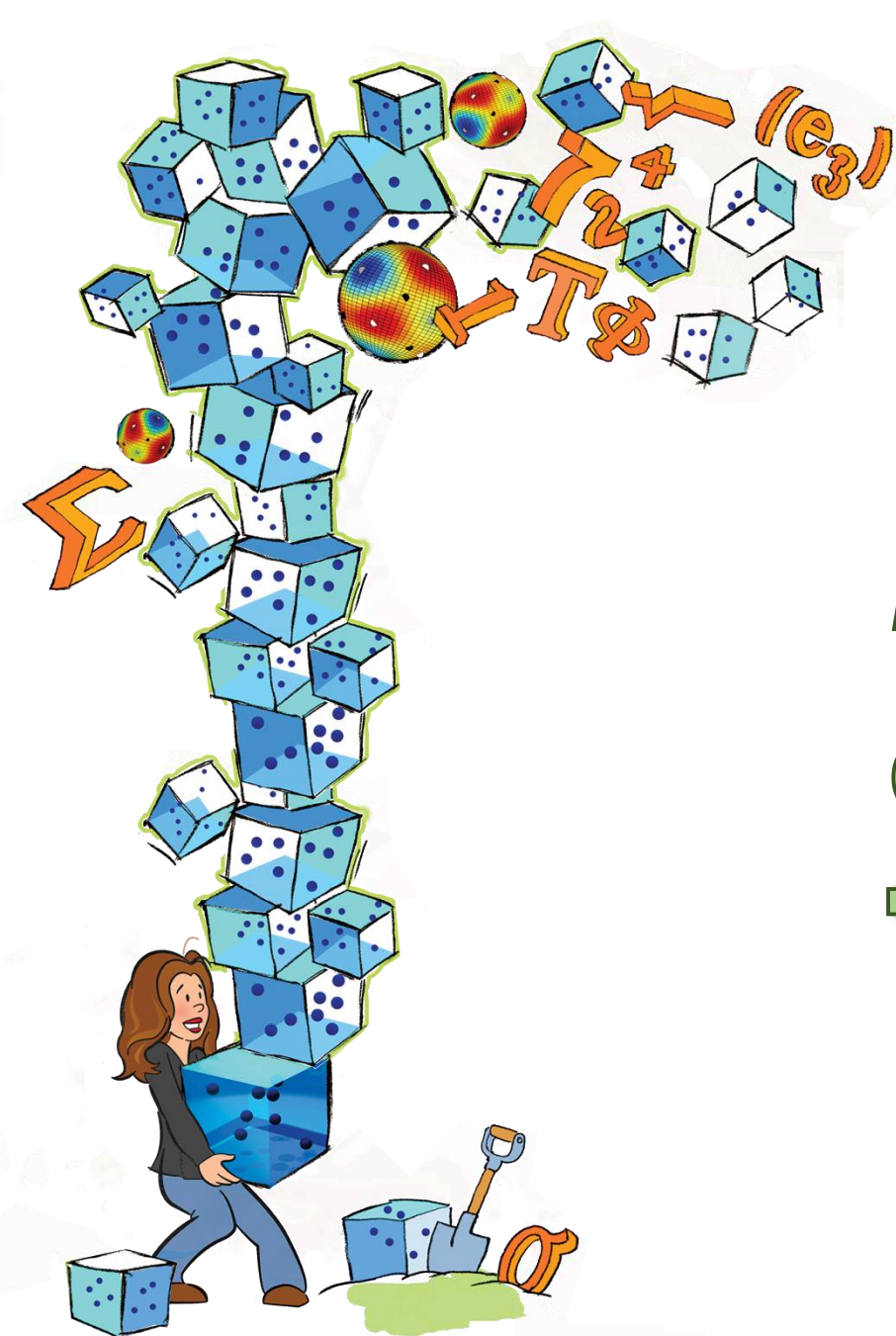
$$\hat{\mathcal{M}}^{(d)} = \frac{1}{p} \sum_{i=1}^p \mathbf{x}_i^{\otimes d}$$

$$\left\| \hat{\mathcal{M}}^{(d)} - \tilde{\mathcal{M}}^{(d)} \right\| = \left\| \frac{1}{p} \sum_{i=1}^p \mathbf{x}_i^{\otimes d} - \sum_{j=1}^m \lambda_j \mu_j^{\otimes d} \right\|$$





More Tensors and Symmetric Tensors



More Tensor Products



$$\mathbf{b} \in \mathbb{R}^n \quad \mathbf{C} \in \mathbb{R}^{n \times n}$$

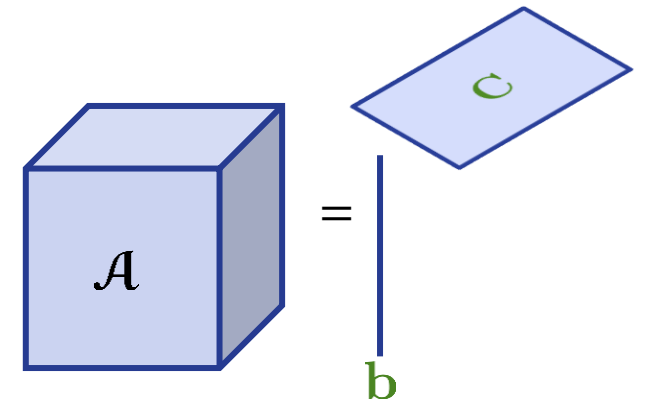
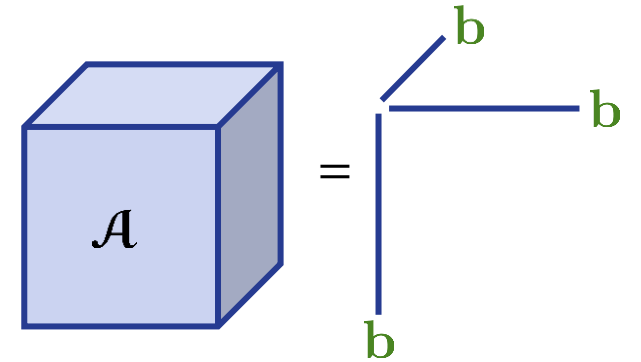
$$\mathcal{A} = \mathbf{b}^{\otimes 3} = \mathbf{b} \otimes \mathbf{b} \otimes \mathbf{b} \in \mathbb{R}^{n \times n \times n} \iff a_{ijk} = b_i b_j b_k$$

$$\mathcal{A} = \mathbf{b} \otimes \mathbf{C} \in \mathbb{R}^{n \times n \times n} \iff a_{ijk} = b_i c_{jk}$$

(not symmetric!)

$$\mathcal{A} = \mathbf{b}^{\otimes 2} \otimes \mathbf{C} \in \mathbb{R}^{n \times n \times n \times n} \iff a_{ijkl} = b_i b_j c_{kl}$$

$$\mathcal{A} = \mathbf{C}^{\otimes 2} = \mathbf{C} \otimes \mathbf{C} \in \mathbb{R}^{n \times n \times n \times n} \iff a_{ijkl} = c_{ij} c_{kl}$$



Symmetrization



$$\text{sym} \left(\begin{array}{c} \text{green top} \\ \text{red front} \\ \text{blue right} \end{array} \right) = \frac{1}{6} \left(\begin{array}{c} \text{green top} \\ \text{red front} \\ \text{blue right} \end{array} + \begin{array}{c} \text{blue top} \\ \text{red front} \\ \text{green right} \end{array} + \begin{array}{c} \text{green top} \\ \text{blue front} \\ \text{red right} \end{array} + \begin{array}{c} \text{blue top} \\ \text{green front} \\ \text{red right} \end{array} + \begin{array}{c} \text{red top} \\ \text{blue front} \\ \text{green right} \end{array} + \begin{array}{c} \text{red top} \\ \text{green front} \\ \text{blue right} \end{array} \right)$$

$$(\text{sym}(\mathcal{A}))_{ijk} = \frac{1}{6} (a_{ijk} + a_{ikj} + a_{jik} + a_{jki} + a_{kij} + a_{kji})$$

Lemma (Hackbusch, 2019)

$$\langle \text{sym}(\mathcal{A}), \mathbf{b}^{\otimes d} \rangle = \langle \mathcal{A}, \mathbf{b}^{\otimes d} \rangle$$

BONUS

Symmetric Tensors Correspond to Polynomials



MathSci.ai

$$\mathcal{A} \in \mathbb{R}^{n \times n \times n}$$

$$\Phi[\mathcal{A}](z_1, \dots, z_n) \equiv \sum_{i=1}^n \sum_{j=1}^n \sum_{k=1}^n a_{ijk} z_i z_j z_k = \langle \mathcal{A}, \mathbf{z}^{\otimes 3} \rangle$$

Proposition (Kileel-K-Pereira 2022)

$$\Phi[\mathbf{a}](z_1, \dots, z_n) = \mathbf{a}^\top \mathbf{z}$$

$$\Phi[\mathbf{A}](z_1, \dots, z_n) = \mathbf{z}^\top \mathbf{A} \mathbf{z}$$

$$\Phi[\mathcal{A}](z_1, \dots, z_n) = \Phi[\text{sym}(\mathcal{A})](z_1, \dots, z_n)$$

$$\Phi[\mathcal{A} \otimes \mathcal{B}](z_1, \dots, z_n) = \Phi[\mathcal{A}](z_1, \dots, z_n) \cdot \Phi[\mathcal{B}](z_1, \dots, z_n)$$

Binomial Theorem for Tensors

(Kileel-K-Pereira 2022)

$$(\mathbf{a} + \mathbf{b})^{\otimes d} = \sum_{k=0}^d \binom{d}{k} \text{sym}(\mathbf{a}^{\otimes k} \otimes \mathbf{b}^{\otimes d-k})$$



Method of Moments for Gaussian Mixture Models – General Scenario

Gaussian Mixture Model: General Case

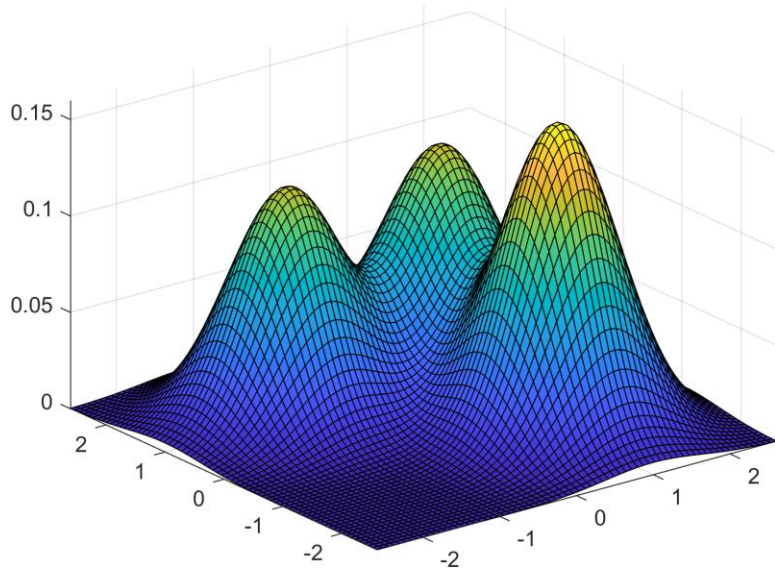


$$X \sim \sum_{j=1}^m \lambda_j \mathcal{N}(\mu_j, \Sigma_j)$$

$$\min_{\theta} f(\theta) \equiv \left\| \frac{1}{p} \sum_{i=1}^p \mathbf{x}_i^{\otimes d} - \mathcal{M}^{(d)} \right\|^2$$

Optimization Parameters

$$\theta = \{ \lambda_j, \mu_j, \Sigma_j \}_{j=1}^m$$



empirical moment tensor

$$\frac{1}{p} \sum_{i=1}^p \mathbf{x}_i^{\otimes d}$$

model moment tensor

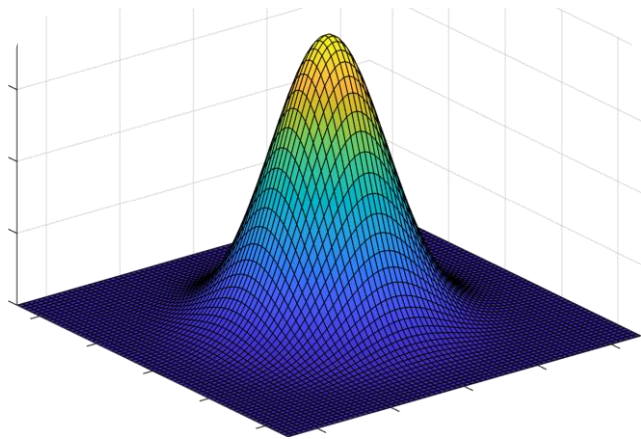
$$\mathcal{M}^{(d)}$$

Problem: Can we explicitly characterize the model moment tensor in terms of the parameters?

Gaussian Model Moment



Let $X \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$, $\mathcal{M}^{(d)} \equiv \mathbb{E}(X^{\otimes d})$. Then



$$\mathcal{M}^{(1)} = \boldsymbol{\mu}$$

$$\mathcal{M}^{(2)} = \boldsymbol{\mu}^{\otimes 2} + \boldsymbol{\Sigma}$$

$$\mathcal{M}^{(3)} = \boldsymbol{\mu}^{\otimes 3} + 3 \operatorname{sym}(\boldsymbol{\mu} \otimes \boldsymbol{\Sigma})$$

$$\mathcal{M}^{(4)} = \boldsymbol{\mu}^{\otimes 4} + 6 \operatorname{sym}(\boldsymbol{\mu}^{\otimes 2} \otimes \boldsymbol{\Sigma}) + 3 \operatorname{sym}(\boldsymbol{\Sigma}^{\otimes 2})$$

Theorem (Kileel-K-Pereira 2022)

$$\mathcal{M}^{(d)} = \sum_{k=0}^{\lfloor d/2 \rfloor} \binom{d}{2k} \frac{2k!}{k!2^k} \operatorname{sym}(\boldsymbol{\mu}^{\otimes d-2k} \otimes \boldsymbol{\Sigma}^{\otimes k})$$

Proof techniques

- Equivalence of symmetric tensors and polynomials
- Marginals of multivariate Gaussians are univariate Gaussian
- Binomial theorem

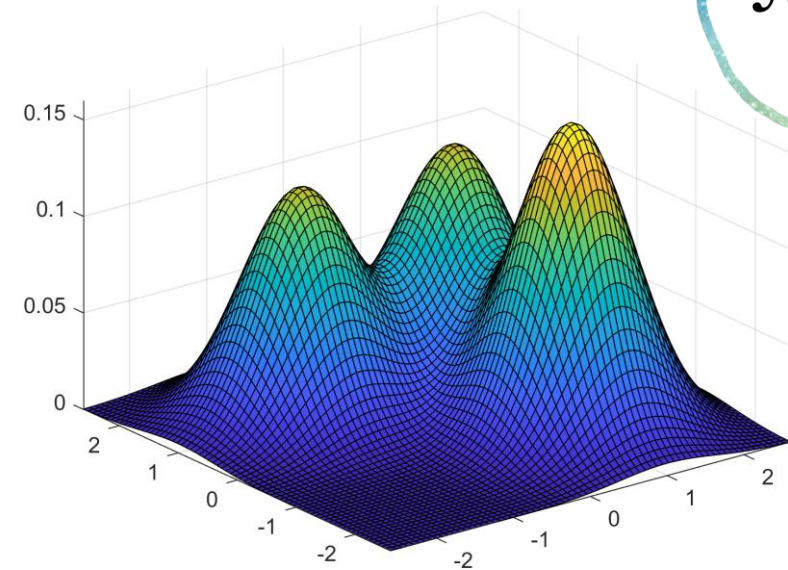
GMM Model Moment

Theorem (Kileel-K-Pereira 2022)

Let $X \sim \sum_{j=1}^m \lambda_j \mathcal{N}(\mu_j, \Sigma_j)$, $\mathcal{M}^{(d)} \equiv \mathbb{E}(X^{\otimes d})$.

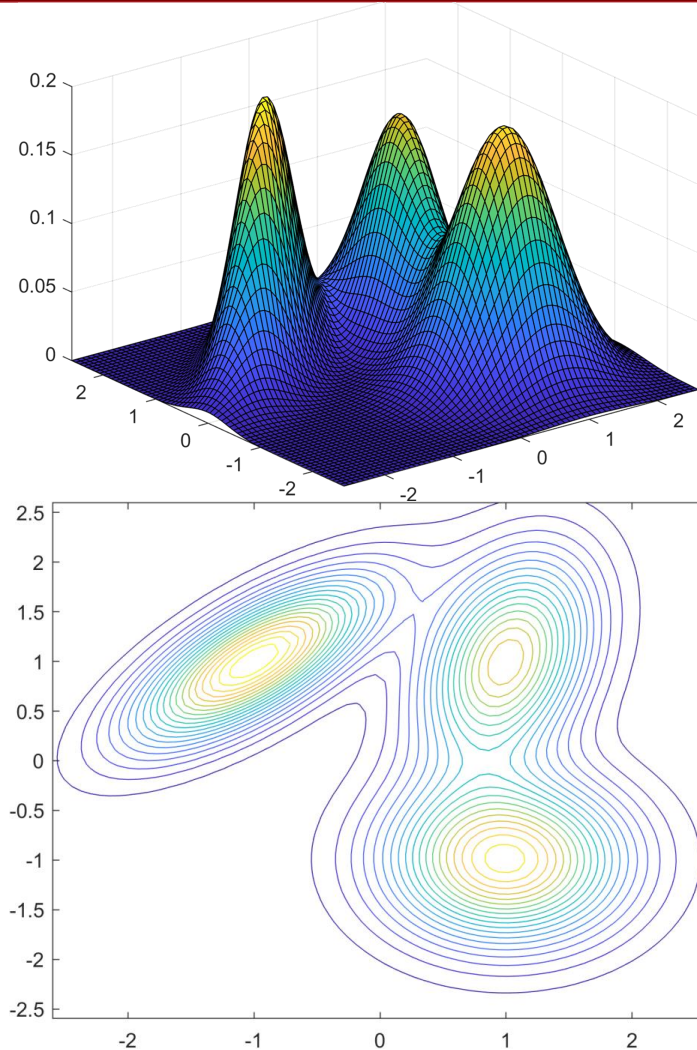
Then

$$\mathcal{M}^{(d)} = \sum_{j=1}^m \lambda_j \mathcal{M}_j^{(d)}$$

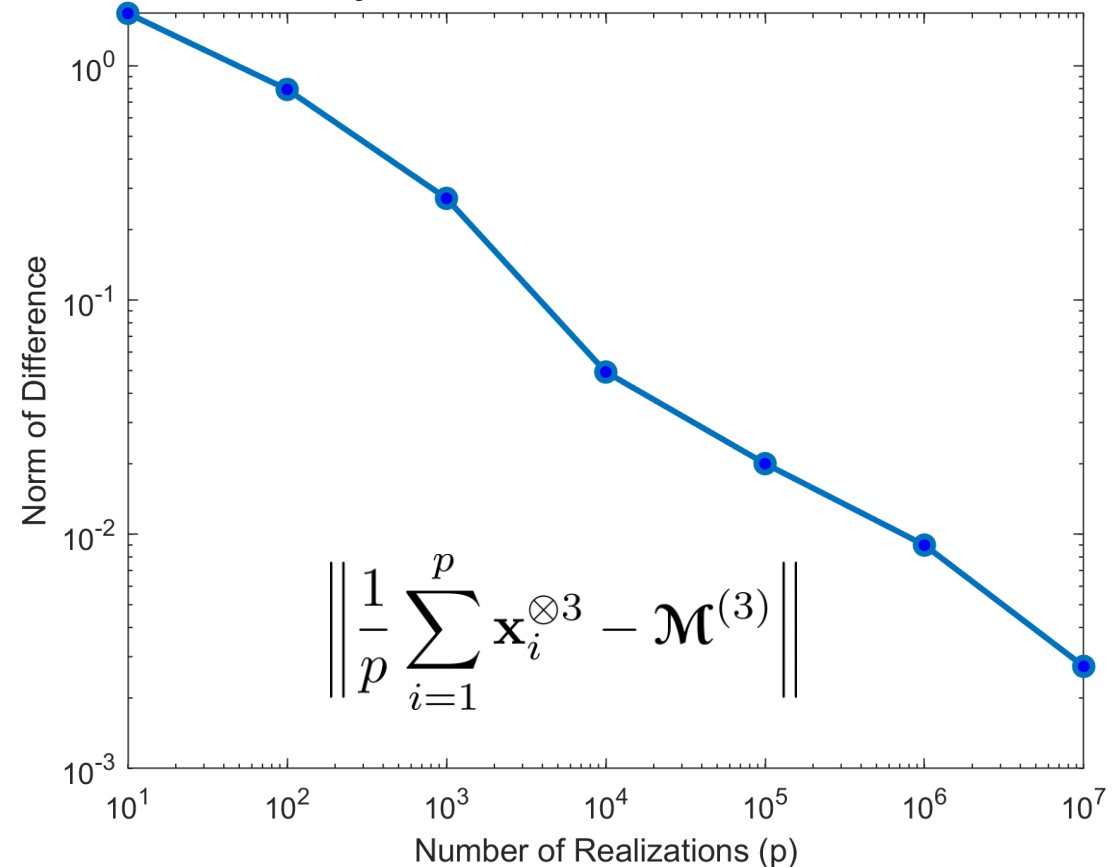


$$\mathcal{M}_j^{(d)} = \sum_{k=0}^{\lfloor d/2 \rfloor} \binom{d}{2k} \frac{2k!}{k! 2^k} \text{sym} \left(\mu_j^{\otimes d-2k} \otimes \Sigma_j^{\otimes k} \right)$$

Empirical versus GMM Model Moment



$$\mathcal{M}^{(3)} = \sum_{j=1}^3 \lambda_j \left(\mu_j^{\otimes 3} + 3 \text{sym}(\mu_j \otimes \Sigma_j) \right)$$



Previous Approach was Biased



New Method: Exact

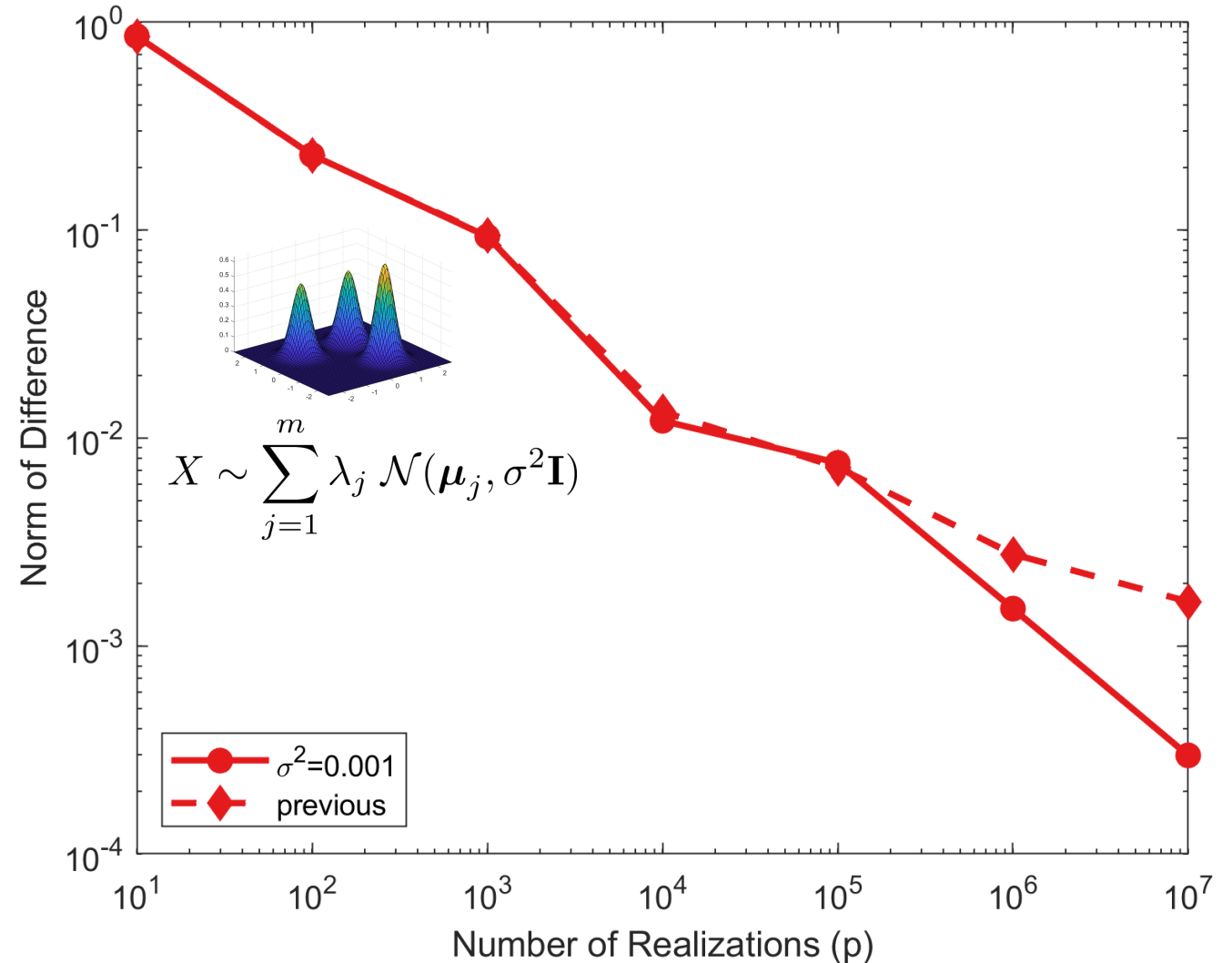
$$\mathcal{M}_j^{(3)} = \mu_j^{\otimes 3} + 3 \text{sym}(\mu_j \otimes \Sigma_j)$$

$$\left\| \frac{1}{p} \sum_{i=1}^p \mathbf{x}_i^{\otimes 3} - \sum_{j=1}^m \lambda_j \mathcal{M}_j^{(3)} \right\|$$

Previous Method: Approximate

$$\widetilde{\mathcal{M}}_j^{(3)} = \mu_j^{\otimes 3}$$

$$\left\| \frac{1}{p} \sum_{i=1}^p \mathbf{x}_i^{\otimes 3} - \sum_{j=1}^m \lambda_j \mu_j^{\otimes 3} \right\|$$



Previous Approach was Biased



New Method: Exact

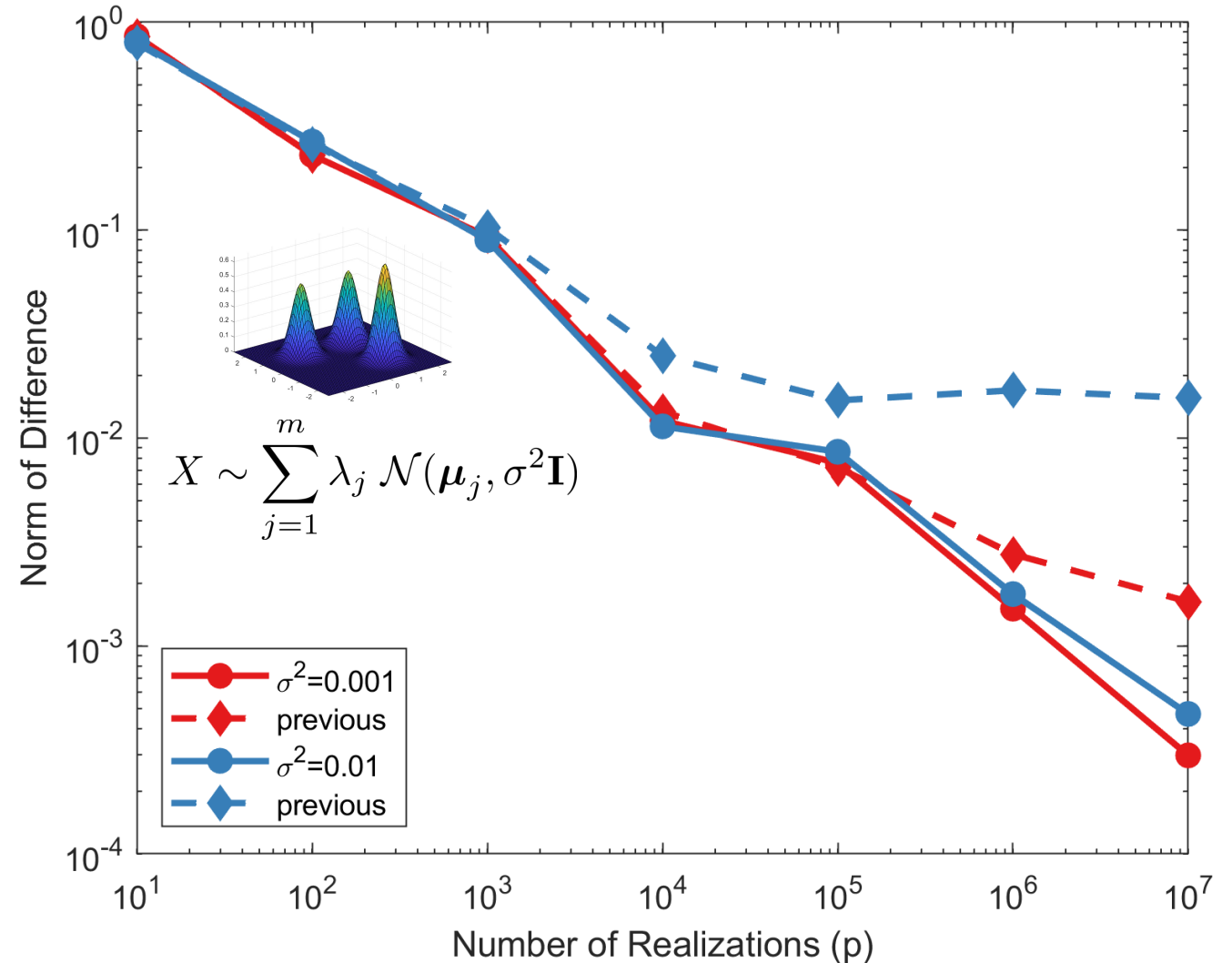
$$\mathcal{M}_j^{(3)} = \mu_j^{\otimes 3} + 3 \text{sym}(\mu_j \otimes \Sigma_j)$$

$$\left\| \frac{1}{p} \sum_{i=1}^p \mathbf{x}_i^{\otimes 3} - \sum_{j=1}^m \lambda_j \mathcal{M}_j^{(3)} \right\|$$

Previous Method: Approximate

$$\widetilde{\mathcal{M}}_j^{(3)} = \mu_j^{\otimes 3}$$

$$\left\| \frac{1}{p} \sum_{i=1}^p \mathbf{x}_i^{\otimes 3} - \sum_{j=1}^m \lambda_j \mu_j^{\otimes 3} \right\|$$



Previous Approach was Biased



New Method: Exact

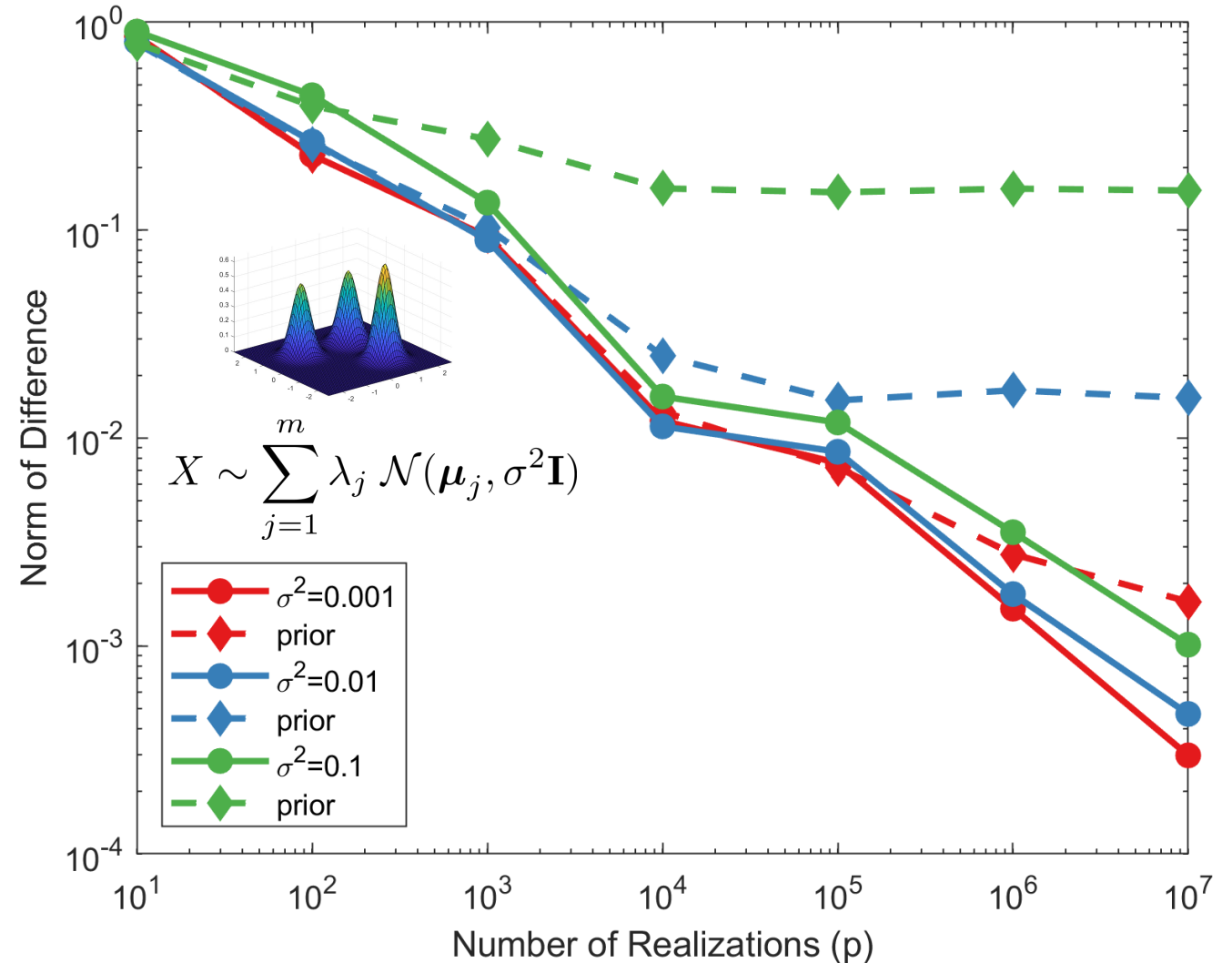
$$\mathcal{M}_j^{(3)} = \boldsymbol{\mu}_j^{\otimes 3} + 3 \text{sym}(\boldsymbol{\mu}_j \otimes \boldsymbol{\Sigma}_j)$$

$$\left\| \frac{1}{p} \sum_{i=1}^p \mathbf{x}_i^{\otimes 3} - \sum_{j=1}^m \lambda_j \mathcal{M}_j^{(3)} \right\|$$

Previous Method: Approximate

$$\widetilde{\mathcal{M}}_j^{(3)} = \boldsymbol{\mu}_j^{\otimes 3}$$

$$\left\| \frac{1}{p} \sum_{i=1}^p \mathbf{x}_i^{\otimes 3} - \sum_{j=1}^m \lambda_j \boldsymbol{\mu}_j^{\otimes 3} \right\|$$



Optimization Formulation



$$X \sim \sum_{j=1}^m \lambda_j \mathcal{N}(\mu_j, \Sigma_j)$$

$$\min_{\theta} f(\theta) \equiv \left\| \frac{1}{p} \sum_{i=1}^p \mathbf{x}_i^{\otimes d} - \mathcal{M}^{(d)} \right\|^2$$

empirical moment tensor

model moment tensor

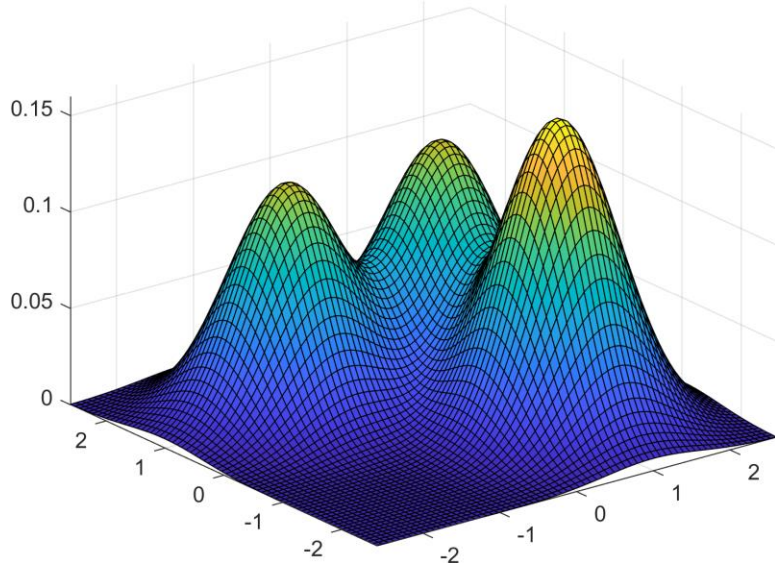
$$\frac{1}{p} \sum_{i=1}^p \mathbf{x}_i^{\otimes d}$$

$$\sum_{j=1}^m \lambda_j \mathcal{M}_j^{(d)}$$

Optimization Parameters

$$\theta = \{ \lambda_j, \mu_j, \Sigma_j \}_{j=1}^m$$

Problem: Can we explicitly characterize the model moment tensor in terms of the parameters?



Yes!

$$\mathcal{M}_j^{(d)} = \sum_{k=0}^{\lfloor d/2 \rfloor} \binom{d}{2k} \frac{2k!}{k!2^k} \text{sym}(\mu_j^{\otimes d-2k} \otimes \Sigma_j^{\otimes k})$$

Optimization Formulation



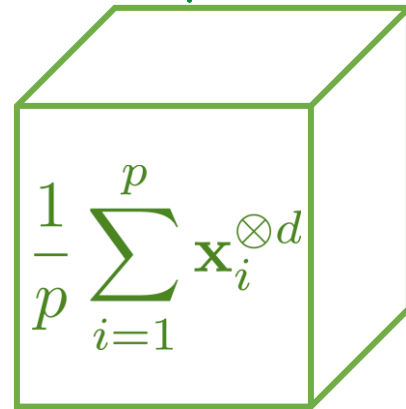
Optimization Parameters

$$\theta = \{ \lambda_j, \mu_j, \Sigma_j \}_{j=1}^m$$

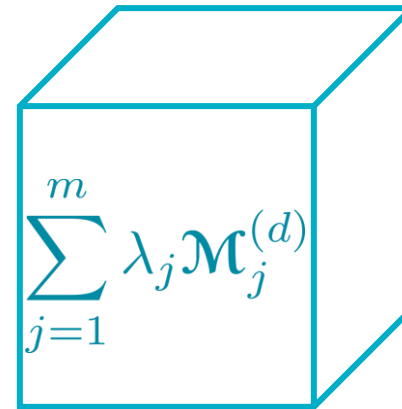
$$\min_{\theta} f(\theta) \equiv \left\| \frac{1}{p} \sum_{i=1}^p \mathbf{x}_i^{\otimes d} - \sum_{j=1}^m \lambda_j \mathcal{M}_j^{(d)} \right\|^2$$



empirical moment tensor



model moment tensor



$$\mathcal{M}_j^{(d)} = \sum_{k=0}^{\lfloor d/2 \rfloor} \binom{d}{2k} \frac{2k!}{k!2^k} \text{sym} \left(\mu_j^{\otimes d-2k} \otimes \Sigma_j^{\otimes k} \right)$$

Optimization Formulation as Inner Products



$$\min_{\theta} f(\theta) \equiv \left\| \frac{1}{p} \sum_{i=1}^p \mathbf{x}_i^{\otimes d} - \sum_{j=1}^m \lambda_j \mathcal{M}_j^{(d)} \right\|^2$$

$$f(\theta) = \left\| \frac{1}{p} \sum_{i=1}^p \mathbf{x}_i^{\otimes d} \right\|^2 + \left\| \sum_{j=1}^m \lambda_j \mathcal{M}_j^{(d)} \right\|^2 - 2 \left\langle \frac{1}{p} \sum_{i=1}^p \mathbf{x}_i^{\otimes d}, \sum_{j=1}^m \lambda_j \mathcal{M}_j^{(d)} \right\rangle$$

constant

$$f(\theta) = C + \sum_{i=1}^m \sum_{j=1}^m \lambda_i \lambda_j \langle \mathcal{M}_i^{(d)}, \mathcal{M}_j^{(d)} \rangle - \frac{2}{p} \sum_{i=1}^p \sum_{j=1}^m \lambda_j \langle \mathbf{x}_i^{\otimes d}, \mathcal{M}_j^{(d)} \rangle$$

dot product of 2 moments

dot product of moment + vector

See forthcoming arXiv posting for full details!!



Example Calculation: $d = 3$

$$f(\theta) = C + \sum_{i=1}^m \sum_{j=1}^m \lambda_i \lambda_j \langle \mathcal{M}_i^{(d)}, \mathcal{M}_j^{(d)} \rangle - \frac{2}{p} \sum_{i=1}^p \sum_{j=1}^m \lambda_j \langle \mathbf{x}_i^{\otimes d}, \mathcal{M}_j^{(d)} \rangle$$

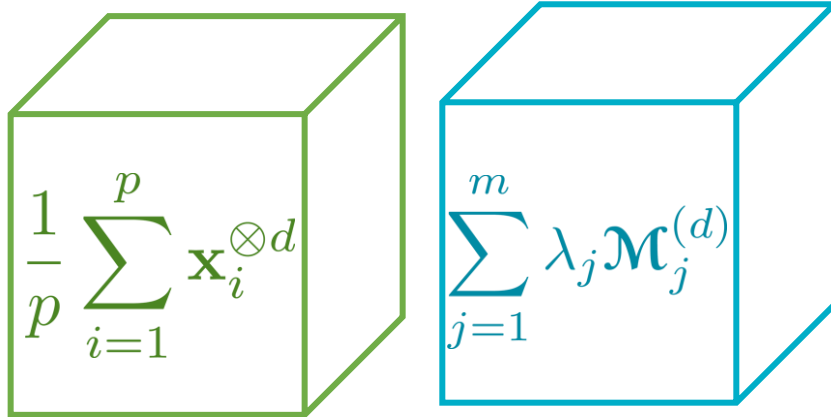
$$\mathcal{M}_j^{(3)} = \boldsymbol{\mu}_j^{\otimes 3} + 3 \operatorname{sym}(\boldsymbol{\mu}_j \otimes \boldsymbol{\Sigma}_j)$$

$$\begin{aligned} \langle \mathbf{x}_i^{\otimes 3}, \mathcal{M}_j^{(3)} \rangle &= \langle \mathbf{x}_i^{\otimes 3}, \boldsymbol{\mu}_j^{\otimes 3} \rangle + 3 \langle \mathbf{x}_i^{\otimes 3}, \operatorname{sym}(\boldsymbol{\mu}_j \otimes \boldsymbol{\Sigma}_j) \rangle \\ &= (\mathbf{x}_i^\top \boldsymbol{\mu}_j)^3 + 3 \langle \mathbf{x}_i^{\otimes 3}, \boldsymbol{\mu}_j \otimes \boldsymbol{\Sigma}_j \rangle \\ &= (\mathbf{x}_i^\top \boldsymbol{\mu}_j)^3 + 3(\mathbf{x}_i^\top \boldsymbol{\mu}_j)(\mathbf{x}_i^\top \boldsymbol{\Sigma}_j \mathbf{x}_i) \end{aligned}$$

Theory for computing inner product of two GMM moment tensors or empirical and GMM moment tensor – see forthcoming arXiv paper for details!



Casting as Optimization Problem



$$\min_{\theta} f(\theta) = C + \sum_{i=1}^m \sum_{j=1}^m \lambda_i \lambda_j \langle \mathcal{M}_i^{(d)}, \mathcal{M}_j^{(d)} \rangle - \frac{2}{p} \sum_{i=1}^p \sum_{j=1}^m \lambda_j \langle \mathbf{x}_i^{\otimes d}, \mathcal{M}_j^{(d)} \rangle$$

$$\theta = \{ \lambda_j, \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j \}_{j=1}^m$$

- Never need to form empirical or model moments explicitly, overcoming curse of dimensionality
- Function can be calculated using simple calculations, total work $\mathcal{O}(m^2n + pmn^2 + m^2n^3)$ per iteration and $\mathcal{O}(mn + pn)$ storage
- Gradients can be calculated as well
- Easy stochastic function and gradient if number of samples (p) is large, as before
- *Issue:* Inherent scaling problem

Non-uniqueness problem



$$X \sim 0.5 \mathcal{N}(\boldsymbol{\mu}_1, \boldsymbol{\Sigma}_1) + 0.5 \mathcal{N}(\boldsymbol{\mu}_2, \boldsymbol{\Sigma}_2)$$

$$Y \sim 0.25 \mathcal{N}\left(\left(\frac{1}{2}\right)^{1/3} \boldsymbol{\mu}_1, \left(\frac{1}{2}\right)^{2/3} \boldsymbol{\Sigma}_1\right) + 0.75 \mathcal{N}\left(\left(\frac{3}{2}\right)^{1/3} \boldsymbol{\mu}_2, \left(\frac{3}{2}\right)^{2/3} \boldsymbol{\Sigma}_2\right)$$

$$\mathbb{E}(X^{\otimes 3}) = \mathbb{E}(Y^{\otimes 3})$$



Non-uniqueness problem

$$X \sim \sum_{j=1}^m \lambda_j \mathcal{N}(\boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j) \quad Y \sim \sum_{j=1}^m \lambda_j \gamma_j^{-d} \mathcal{N}(\gamma_j \boldsymbol{\mu}_j, \gamma_j^2 \boldsymbol{\Sigma}_j)$$

Problem: Moments of order d are the same!

FIX: Append a constant c to the end of every observation vector, creating vectors of dimension $n + 1$

RESULT: *Implicitly*, a weighted combination of *all the moments* from 1 to d . This means we include all moments up to order d in the optimization.

See forthcoming arXiv posting for full details!!

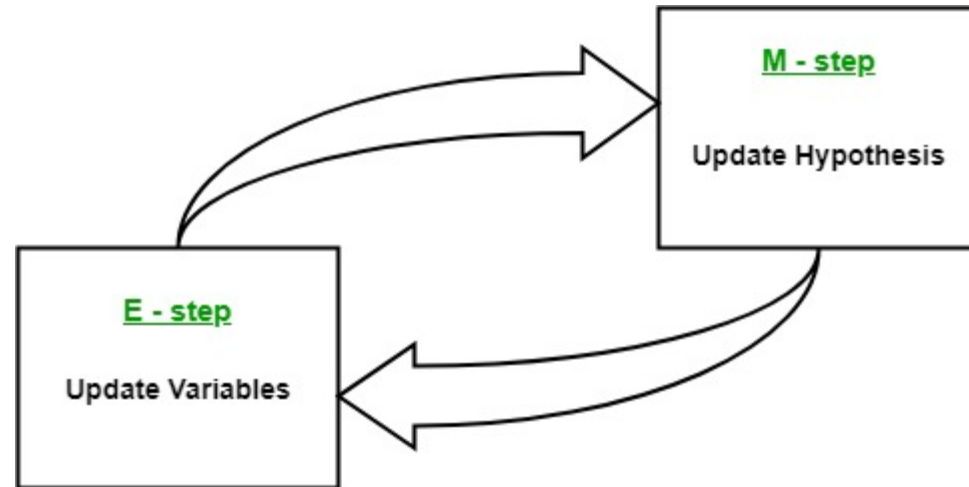


Numerical Results and Conclusions





State of the Art: Expectation Maximization (EM)



EM is State of the Art

- Inexpensive
- Relatively easy to implement
- Optimizing a different cost function
- Sensitive to initialization
- Sensitive to overlapping Gaussians

MoM has theoretical advantages but has not been used much in practice previously because of its great expense

EM Algorithm in a Nutshell

Make initial guesses for parameters

Repeat until log-likelihood converges:

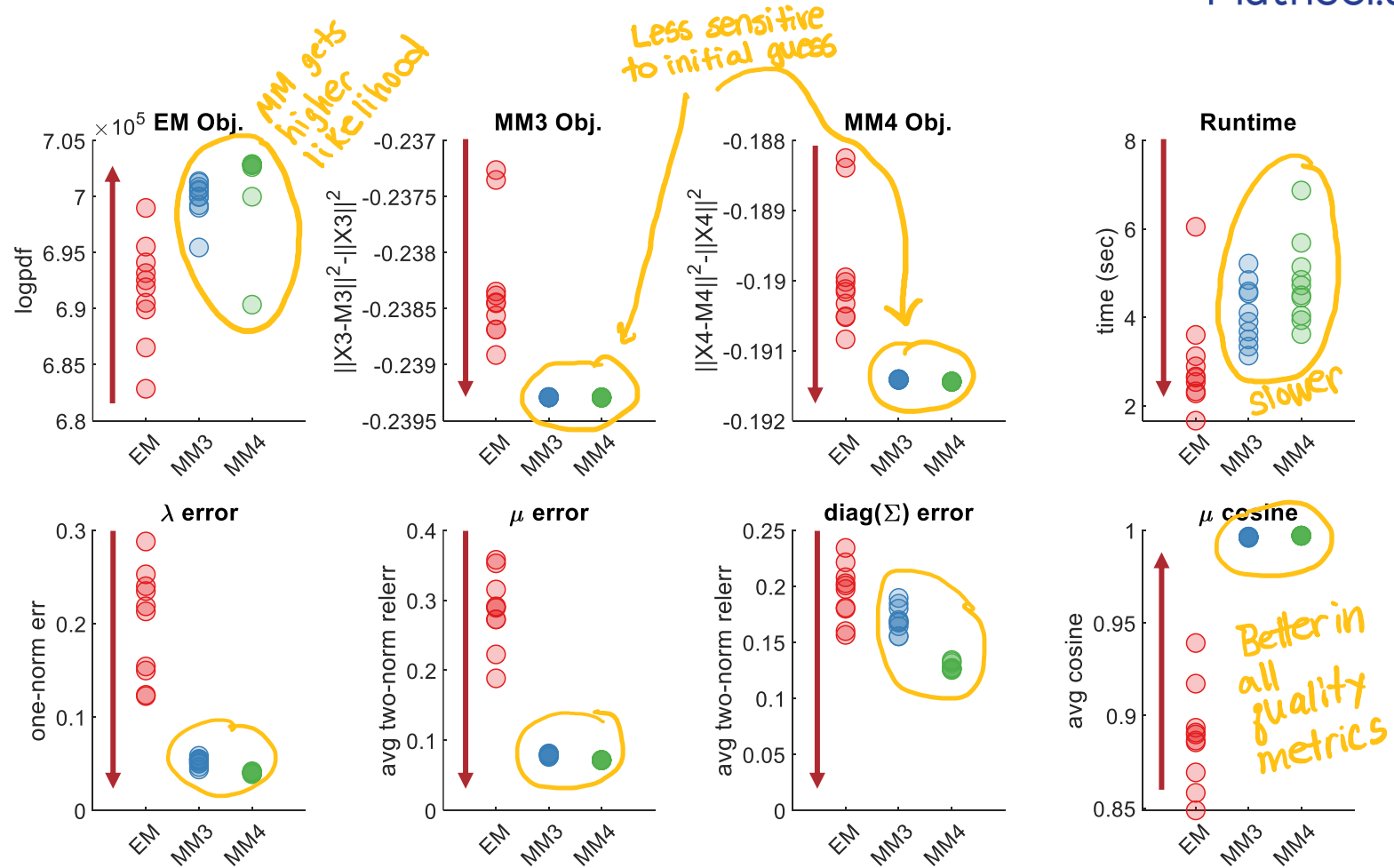
1. Compute membership weights for each datapoint
2. Update the component parameters using the membership weights

See, e.g., *Xu and Jordan (1996)* for discussion of its robustness



Method of Moments can beat EM

- Randomly-generated problem with overlapping Gaussians
 - Diagonal covariances
 - Dimensionality: $n = 100$
 - Number of Gaussians: $m = 20$
 - Observations: $p = 8000$
- Compared three methods
 - EM: Expectation Maximization
 - MM3: Method of Moments, $d = 3$
 - MM4: Method of Moments, $d = 4$
- 10 runs each with different initial guesses

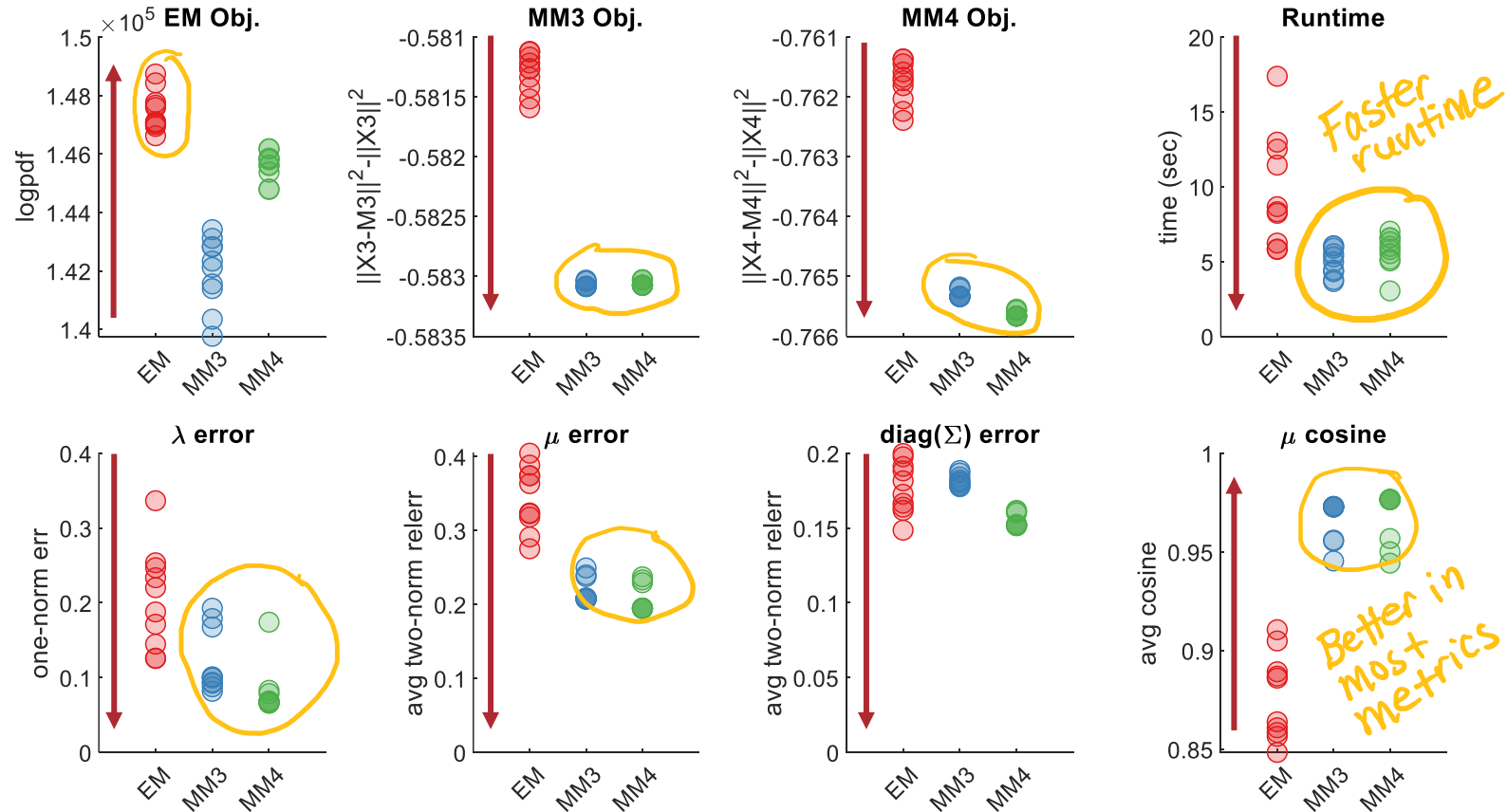




Method of Moments can beat EM

Same setup as previous slide except higher noise

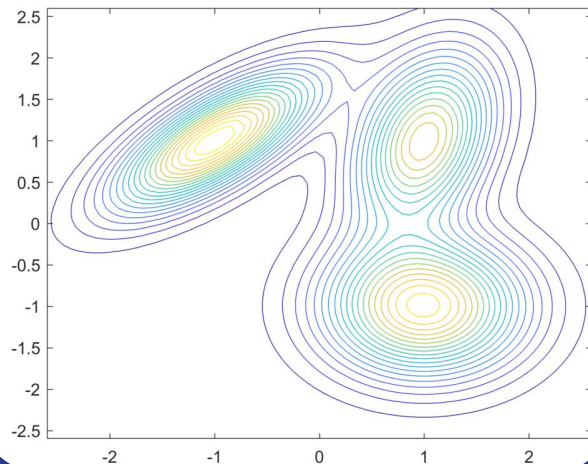
- Randomly-generated problem with overlapping Gaussians
 - Diagonal covariances
 - Dimensionality: $n = 100$
 - Number of Gaussians: $m = 20$
 - Observations: $p = 8000$
- Compared three methods
 - EM: Expectation Maximization
 - MM3: Method of Moments, $d = 3$
 - MM4: Method of Moments, $d = 4$
- 10 runs each with different initial guesses



Related Works

Key Differences in Our Work

- Novel tensor formulation of Gaussian moments
- No spherical or axis-aligned covariance assumptions
- Computationally efficient, no exponential dependence on d



- Many theoretical advantages to method of moments and connections to tensors
 - Hsu and Kakade (2013) – diagonal covariance, $d \leq 3$
 - Ge, Huang, Kakade (2015) – vectorized covariances, loses symmetries
 - Bakshi, Diakonikolas, Jia, Kane, Kothari, and Vempala (2020) – robust learning using tensor decomposition
 - Khouga, Mattei, Mourrian (2021) – GMM identifiability
- Computational approaches (limited handling of covariances)
 - Anandkumar, Ge, Hsu, Kakade (2014) and Anandkumar, Ge, Hsu, Kakade, Telegarsky (2014) – orthogonal symmetric tensor decomposition
 - Sherman & K., 2020 – (general) symmetric tensor decomposition, emphasis of implicit computation to avoid curse of dimensionality
- Inner products of moment tensors
 - Muandet, Fukumizu, Dinuzzo, Schölkopf (2012) – up to 3rd order



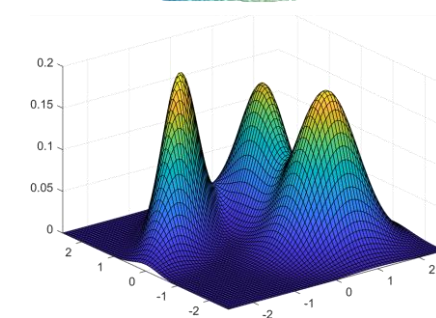
Take-aways and Future Work

- Our focus: Method of moments for Gaussian mixture models (GMMs)

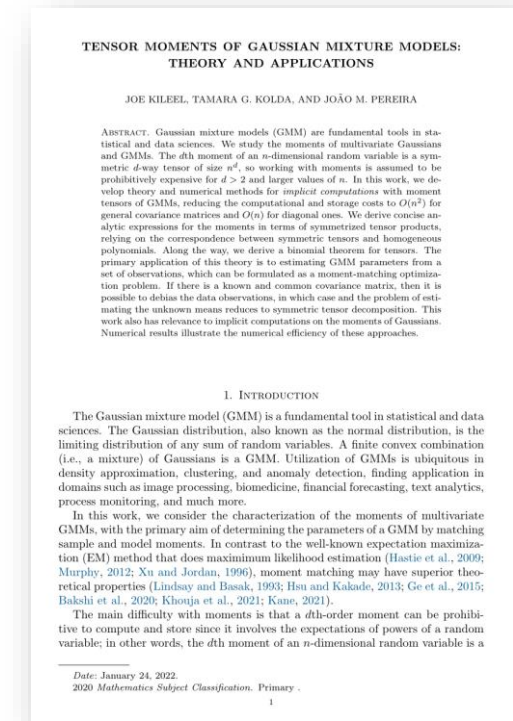
- Key results

- Formulation of GMM moment in terms of tensor outer products

$$\mathcal{M}^{(d)} = \sum_{j=1}^m \sum_{k=0}^{\lfloor d/2 \rfloor} \binom{d}{2k} \frac{2k!}{k!2^k} \lambda_j \text{sym} \left(\mu_j^{\otimes d-2k} \otimes \Sigma_j^{\otimes k} \right)$$



See forthcoming arXiv posting for full details!!



- Efficient computation and storage, avoiding exponential dependence on moment order
 - Novel approach to scaling ambiguity using augmentation
 - Amenable to stochastic formulations
 - Plus...dot product of moment tensors in terms of Bell polynomials, avoiding exponential dependence on moment order
 - Plus...modifying empirical moment tensor to “remove” Gaussian noise

- Future work

- Implementation details, especially for general Gaussians
 - Analysis of optimization landscape and comparison to that of max likelihood
 - Bounding number of samples required for accurate estimation
 - Application studies

J. Kileel, T. Kolda, and J. M. Pereira.
Tensor Moments of Gaussian Mixture Models: Theory and Applications