Singularity structures and parameter estimation in mixture models

Long Nguyen

Department of Statistics University of Michigan

Institute for Mathematical Sciences National University of Singapore June 2016

Joint work with Nhat Ho (UM)

(4) (3) (4) (4) (4)

- machine learning: "model-free optimization-based algorithms"
 - isn't it the spirit of empirical likelihood based methods?
 - prediction vs estimation/inference
- Bayesian nonparametrics:
 - how to be Bayesian, yet more empirical by being nonparametric

(4) (3) (4) (4) (4)

- machine learning: "model-free optimization-based algorithms"
 - isn't it the spirit of empirical likelihood based methods?
 - prediction vs estimation/inference
- Bayesian nonparametrics:
 - how to be Bayesian, yet more empirical by being nonparametric
- empirical likelihood:
 - how to take the inference a bit beyond empirical distributions

(4) (1) (4) (4) (4)

- machine learning: "model-free optimization-based algorithms"
 - isn't it the spirit of empirical likelihood based methods?
 - prediction vs estimation/inference
- Bayesian nonparametrics:
 - how to be Bayesian, yet more empirical by being nonparametric
- empirical likelihood:
 - how to take the inference a bit beyond empirical distributions

modern statistics/ data science

- data increasingly high-dimensional and complex
- inferential goals increasingly more ambitious
- requiring more sophisticated algorithms and complex statistical modeling

A (10) × (10) × (10)

Mixture modeling



Mixture density

$$p(x) = \sum_{i=1}^{k} p_i f(x|\eta_i)$$

e.g., $f(x|\eta_i) = \text{Normal}(x|\mu_i, \Sigma_i)$

< □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ >

Parameter estimation

mixing probabilities $\mathbf{p} = (p_1, \dots, p_k)$? mixing components $\boldsymbol{\eta} = (\eta_1, \dots, \eta_k)$?

Hierarchical models

 $\mathsf{HM} = \mathsf{Mixture} \text{ of mixture models}$

Challenge: parameter estimation in latent variable models



[courtesy M. Jordan's slides]

< □ > < 同 > < 回 > < 回 > < 回 >

Estimation in parametric models

Standard methods, e.g., maximum likelihood estimation (via EM) or Bayesian estimation, yield root-*n* parameter estimation rate

A (10) × A (10) × A (10)

Estimation in parametric models

Standard methods, e.g., maximum likelihood estimation (via EM) or Bayesian estimation, yield root-n parameter estimation rate

• if Fisher information matrix is non-singular

What if we are in a singular situation?

Fisher singularity

- Cox & Hinkley (1974), Lee & Chesher (1986), Rotnitzky, Cox, Bottai and Robbins (2000), etc: first-order singularity
- Azzalini & Capitanio (1999); Pewsey (2000); DiCiccio & Monti (2004); Hallin & Ley (2012,2014), etc: third order singularities in skewnormal distributions
- Chen (1995); Rousseau & Mengersen (2011); Nguyen (2013): asymptotics for parameter estimation in overfitted mixture models, under strong identifiability conditions
- A full picture of singularity structures in mixture models remain largely unknown (e.g., hitherto there's no asymptotic theory for finite mixtures of location-scale Gaussian mixtures)

イロト 不得 トイラト イラト 一日

Singularity is a common occurence in modern statistics

- high-dimensional and sparse setting
- "complicated" density classes (e.g., gamma, normal, skewnormal), when there are more than one parameter varying
- overfitted/infinite mixture models
- hierarchical models

() < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < ()

Singularities in e-mixtures

▲ロト ▲圖ト ▲国ト ▲国

e-mixtures = exact-fitted mixtures

Mixture model indexed by mixing distribution $G = \sum_{i=1}^{k} p_i \delta_{\eta_i}$

$$\bigg\{p_G(x) = \sum_{i=1}^k p_i f(x|\eta_i) \bigg| G \text{ has } k \text{ atoms} \bigg\},\$$

Fisher information matrix

$$I(G) = \mathbb{E}\left\{\left(\frac{\partial \log p_G(X)}{\partial G}\right) \left(\frac{\partial \log p_G(X)}{\partial G}\right)^T\right\}$$

where $\partial \log p_G / \partial G$ simply denotes partial derivative of score function wrt all parameters **p**, η

< □ > < 同 > < 回 > < 回 > < 回 >

Exploiting the representation $p_G = \sum p_i f(x|\eta_i)$, easy to note that I(G) is non-singular iff the collection of

$$\left\{f(x|\eta_i), rac{\partial f}{\partial \eta}(x|\eta_i)\middle|i=1,\ldots,k
ight\}$$

are linearly independent functions of x

This condition holds if

- $f(x|\eta) = \text{Gaussian}(x|\mu, \sigma^2)$ location-scale Gaussian
- not for Gamma or skewnormal kernels (and many others)

< □ > < 同 > < 回 > < 回 > < 回 >

Gamma mixtures

Gamma density

$$f(x|a,b) = \frac{b^a}{\Gamma(a)} x^{a-1} e^{-bx}, a > 0, b > 0$$

admits the following pde

$$\frac{\partial f}{\partial b}(x|a,b) = \frac{a}{b}f(x|a,b) - \frac{a}{b}f(x|a+1,b).$$

For Gamma mixture model

$$p_G(x) = \sum_{i=1}^k p_i f(x|a_i, b_i)$$

I(G) is a singular matrix if $a_i - a_j = 1$ and $b_i = b_j$ for some pair of components i, j = 1, ..., k.

< □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > <

Theorem (Ho & Nguyen, 2016) The minimax rate of estimation for Gamma mixture's parameter is slower than $\frac{1}{n^{1/2r}}$ for any $r \ge 1$, sample size n

This is because we can't tell very well when G is singular or not

However if we know that the true G is non-singular, and that $|a_i - a_j| - 1$ and $|b_i - b_j|$ is bounded away from 0, then MLE achieves root-n rate

イロト イポト イヨト イヨト 二日

Parameter estimation rate for Gamma mixtures



L: MLE restricted to compact set of non-singular G: $W_1(\widehat{G}_n, G_0) \simeq n^{-1/2}$. R: MLE for general set of G: $W_1 \approx 1/(\log n)^{1/2}$.

Skewnormal mixtures

Skewnormal kernel density

$$f(x|\theta,\sigma,m) := \frac{2}{\sigma}f\left(\frac{x-\theta}{\sigma}\right)\Phi(m(x-\theta)/\sigma),$$

where f(x) is the standard normal density and $\Phi(x) = \int f(t) \mathbf{1}(t \le x) dt$. This generalizes Gaussian densities, which correspond to m = 0.

I(G) is singular iff the parameters are real solution of a number of polynomial equations

(i) Type A:
$$P_1(\eta) = \prod_{j=1}^k m_j$$
.

(ii) Type B:
$$P_2(\eta) = \prod_{1 \le i \ne j \le k} \left\{ (\theta_i - \theta_j)^2 + \left[\sigma_i^2 (1 + m_j^2) - \sigma_j^2 (1 + m_i^2) \right]^2 \right\}$$

・ 何 ト ・ ヨ ト ・ ヨ ト



Figure : Illustration of type A and type B singularity.

< □ > < □ > < □ > < □ > < □ >

- Singularities of skewnormal mixtures lie in affine varieties
- If we know that the parameters are bounded away from these algebraic sets, then method such as MLE continues to produce root-n rate
 Otherwise it is worse, especially if the true model is near singularity
 - in practice, want to test if the true parameters are a singular point
 - in theory, we want to know the actual rate of estimation for singular points, necessitating the need to look into deep structure of singularities

- Singularities of skewnormal mixtures lie in affine varieties
- If we know that the parameters are bounded away from these algebraic sets, then method such as MLE continues to produce root-n rate Otherwise it is worse, especially if the true model is near singularity
 - in practice, want to test if the true parameters are a singular point
 - in theory, we want to know the actual rate of estimation for singular points, necessitating the need to look into deep structure of singularities
- We'll show that the singularity structure is very rich and go beyond the singularities of Fisher information
 - introduce singularity levels, which provide multi-level partitions of parameter space
- There is a consequence: the "more" singular the parameter values, the worse the MLE and minimax rates of estimation, $n^{-1/2}$ to $n^{-1/4}$ to $n^{-1/8}$, ad infinitum

イロト 不得 トイヨト イヨト 二日

General theory

・ロト ・日 ・ ・ ヨト ・

General theory

behavior of likelihood function in the neighborhood of model parameters

< 口 > < 同

→ Ξ →

From parameter space to space of mixing measure G

The map $(\mathbf{p}, \boldsymbol{\eta}) \mapsto G(\mathbf{p}, \boldsymbol{\eta}) = \sum_i p_i \delta_{\eta_i}$ is many-to-one, e.g.

$$\frac{1}{2}\delta_0 + \frac{1}{2}\delta_1 = \frac{1}{2}\delta_0 + \frac{1}{3}\delta_1 + \frac{1}{6}\delta_1$$

Say $({\bf p},\eta)$ and $({\bf p}',\eta')$ are equivalent if the corresponding mixing measures are equal, e.g.,

$$[\mathbf{p}, \theta] = [(1/2, 1/2), (0, 1)] \equiv [(1/2, 1/3, 1/6), (0, 1, 1)]$$

< □ > < 同 > < 回 > < 回 > < 回 >

$$G = \sum p_i \delta_{\boldsymbol{\eta}_i} \mapsto p_G(\cdot) = \sum p_i f(\cdot | \boldsymbol{\eta}_i)$$

・ロト・西ト・モン・ビー シック

Wasserstein space of measures/ optimal transport distance

Optimal transportation problem (Monge-Kantorovich)

how to transport good products from a collection of producers to a collection of consumers located in a common space

how to move the mass from one distribution to another?



squares: locations of producers; circles: locations of consumers

Optimal transport/Wasserstein distance is the optimal cost of transportation of mass from — "production distribution" — to — "consumption distribution"

Long Nguyen (UM)

・ロト ・四ト ・ヨト ・ヨト

Wasserstein distance

Let G, G' be two prob. measures on \mathbb{R}^d , $r \ge 1$. A coupling κ of G, G' is a joint dist on $\mathbb{R}^d \times \mathbb{R}^d$ which induces marginals G, G'. κ also called a "transportation plan".

The *r*-th order Wasserstein distance, denoted by W_r , is given by

$$W_r(G, G') := \left[\inf_{\kappa} \int \|\theta - \theta'\|^r d\kappa(\theta, \theta')\right]^{1/r}.$$

$$G = \sum p_i \delta_{\boldsymbol{\eta}_i} \mapsto p_G(\cdot) = \sum p_i f(\cdot | \boldsymbol{\eta}_i)$$

・ロト・西ト・モン・ビー シック

Behavior of likelihood in a Wasserstein neighborhood

As $G \rightarrow G_0$ in Wasserstein metric, apply Taylor expansion up to the *r*-th order:

$$p_{G}(x) - p_{G_{0}}(x) = \sum_{i=1}^{k_{0}} \sum_{j=1}^{s_{i}} p_{ij} \sum_{|\kappa|=1}^{r} \frac{(\Delta \eta_{ij})^{\kappa}}{\kappa!} \frac{\partial^{|\kappa|} f}{\partial \eta^{\kappa}}(x|\eta_{i}^{0}) + \sum_{i=1}^{k_{0}} \Delta p_{i}f(x|\eta_{i}^{0}) + R_{r}(x),$$

where $R_r(x)$ is the Taylor remainder and $R_r(x)/W_r^r(G,G_0) \rightarrow 0$.

We have seen examples where the partial derivatives up to the first order are not linearly independent (for gamma kernel and skewnormal kernel)

イロト イポト イヨト イヨト 二日

For normal kernel density $f(x|\mu, v)$,

$$\frac{\partial^2 f}{\partial \theta^2} = 2 \frac{\partial f}{\partial v}.$$

< □ > < □ > < □ > < □ > < □ >

For normal kernel density $f(x|\mu, v)$,

$$\frac{\partial^2 f}{\partial \theta^2} = 2 \frac{\partial f}{\partial v}.$$

For skewnormal kernel density $f(x|\mu, v, m)$,

$$\frac{\partial^2 f(x|\eta)}{\partial \theta^2} - 2 \frac{\partial f(x|\eta)}{\partial v} + \frac{m^3 + m}{v} \frac{\partial f(x|\eta)}{\partial m} = 0.$$

< □ > < □ > < □ > < □ > < □ >

r-canonical form

Fix $r \geq 1$. For some sequence of $G_n \in \mathcal{G}$, such that $W_r(G_n, G_0) \to 0$,

$$\frac{p_{G_n}(x) - p_{G_0}(x)}{W_r^r(G_n, G_0)} = \sum_{l=1}^{L_r} \left(\frac{\xi_l(G_0, \Delta G_n)}{W_r^r(G, G_0)} \right) H_l(x|G_0) + o(1),$$

where

- *H*_I(*x*|*G*₀) are *linearly independent* functions
- coefficients ξ_l(G₀, ΔG_n)/W^r_r(G_n, G₀) are ratio of two semi-polynomials of the parameter perturbation of G_n around G₀

< □ > < 同 > < 回 > < Ξ > < Ξ

 H_l may be obtained by reducing partial derivatives to linearly independent ones For Gamma kernel f,

$$\frac{\partial f(x|\eta_j^0)}{\partial m} = -\sum_{j=1}^k \frac{\alpha_{1j}}{\alpha_{4k}} f(x|\eta_j^0) + \frac{\alpha_{2j}}{\alpha_{4k}} \frac{\partial f(x|\eta_j^0)}{\partial \theta} + \frac{\alpha_{3j}}{\alpha_{4k}} \frac{\partial f(x|\eta_j^0)}{\partial v} - \sum_{j=1}^{k-1} \frac{\alpha_{4j}}{\alpha_{4k}} \frac{\partial f(x|\eta_j^0)}{\partial m}$$

• • • • • • • • • • • •

For Gaussian kernel $f(x|\eta) = f(x|\theta, v)$ all partial derivatives wrt both θ and v can be eliminated via the following reduction: for any $\kappa_1, \kappa_2 \in \mathbb{N}$, for any $j = 1, \ldots, k_0$, $\frac{\partial^{\kappa_1 + \kappa_2} f(x|n_i^0)}{\partial^{\kappa_1 + \kappa_2} f(x|n_i^0)} = 1 \quad \frac{\partial^{\kappa_1 + 2\kappa_2} f(x|n_i^0)}{\partial^{\kappa_1 + 2\kappa_2} f(x|n_i^0)}$

$$rac{\partial^{\kappa_1+\kappa_2}f(x|\eta_j^0)}{\partial heta^{\kappa_1}v^{\kappa_2}}=rac{1}{2^{\kappa_2}}rac{\partial^{\kappa_1+2\kappa_2}f(x|\eta_j^0)}{\partial heta^{\kappa_1+2\kappa_2}}.$$

Thus, this reduction is valid for all parameter values (\mathbf{p}, η) , and *r*-canonical forms for all orders.

イロト イポト イヨト イヨト 二日

For skewnormal kernel $f(x|\eta) = f(x|\theta, v, m)$ for any $j = 1, ..., k_0$, any $\eta = \eta_j^0 = (\theta_j^0, v_j^0, m_j^0)$, $\frac{\partial^2 f(x|\eta)}{\partial \theta^2} = 2\frac{\partial f(x|\eta)}{\partial v} - \frac{m^3 + m}{v}\frac{\partial f(x|\eta)}{\partial m}.$

For higher order partial derivatives:

$$\frac{\partial^3 f}{\partial \theta^3} = 2 \frac{\partial^2 f}{\partial \theta \partial v} - \frac{m^3 + m}{v} \frac{\partial^2 f}{\partial \theta \partial m},$$

$$\frac{\partial^3 f}{\partial \theta^2 \partial v} = 2 \frac{\partial^2 f}{\partial v^2} + \frac{m^3 + m}{v^2} \frac{\partial f}{\partial m} - \frac{m^3 + m}{v} \frac{\partial^2 f}{\partial v \partial m},$$

$$\frac{\partial^3 f}{\partial \theta^2 \partial m} = 2 \frac{\partial^2 f}{\partial v \partial m} - \frac{3m^2 + 1}{v} \frac{\partial f}{\partial m} - \frac{m^3 + m}{v} \frac{\partial^2 f}{\partial m^2}.$$

イロト 不得 トイヨト イヨト 二日

r-singularity

Let ${\mathcal G}$ be a space of mixing measure.

Definition

For each finite $r \ge 1$, say G_0 is *r*-singular relative to \mathcal{G} if G_0 admits a *r*-canonical form for some sequence of $G \in \mathcal{G}$, according to which $W_r(G, G_0) \to 0$, coefficients $\xi_l^{(r)}(G)/W_r^r(G, G_0) \to 0$ for all $l = 1, \ldots, L_r$.

< □ > < 同 > < 回 > < 回 > < 回 >

r-singularity

Let ${\mathcal G}$ be a space of mixing measure.

Definition

For each finite $r \ge 1$, say G_0 is *r*-singular relative to \mathcal{G} if G_0 admits a *r*-canonical form for some sequence of $G \in \mathcal{G}$, according to which $W_r(G, G_0) \to 0$, coefficients $\xi_l^{(r)}(G)/W_r^r(G, G_0) \to 0$ for all $l = 1, \ldots, L_r$.

Lemma

(a) The notion of r-singularity is independent of the specific r-form. That is, the existence of the sequence G in the definition holds for all r-canonical forms once it holds for at least one of them.

イロト イポト イヨト イヨト 二日

r-singularity

Let ${\mathcal G}$ be a space of mixing measure.

Definition

For each finite $r \ge 1$, say G_0 is *r*-singular relative to \mathcal{G} if G_0 admits a *r*-canonical form for some sequence of $G \in \mathcal{G}$, according to which $W_r(G, G_0) \to 0$, coefficients $\xi_l^{(r)}(G)/W_r^r(G, G_0) \to 0$ for all $l = 1, \ldots, L_r$.

Lemma

(a) The notion of r-singularity is independent of the specific r-form. That is, the existence of the sequence G in the definition holds for all r-canonical forms once it holds for at least one of them. (b) If G_0 is r-singular for some r > 1, then G_0 is (r - 1)-singular.

イロト 不得 トイラト イラト 一日

Definition

The singularity level of G_0 relative to ambient space \mathcal{G} , denoted by $\ell(G_0|\mathcal{G})$, is

0, if G_0 is not *r*-singular for any $r \ge 1$;

 ∞ , if G_0 is *r*-singular for all $r \ge 1$;

otherwise, the largest natural number $r \ge 1$ for which G_0 is *r*-singular.

< (日) × (日) × (1)

Theorem

If $\ell(G_0|\mathcal{G}) = r$, then for any $G \in \mathcal{G}$ subject to a compactness condition

 $V(p_G, p_{G_0}) \gtrsim W_{r+1}^{r+1}(G, G_0).$

So, singularity level r implies that MLE has rate

 $n^{-\frac{1}{2(r+1)}}$

Under additional regularity condition, this is also a local minimax lower bound for estimating G_0

イロト イポト イヨト イヨト

Role of ambient spaces

If $\mathcal{G}_1 \subset \mathcal{G}_2$ then

$$\ell(G_0|\mathcal{G}_1) \leq \ell(G_0|\mathcal{G}_2)$$

For location-scale Gaussian e-mixtures

 $\ell(G_0|\mathcal{E}_{k_0})=0.$

• • • • • • • • • • • •

Role of ambient spaces

If $\mathcal{G}_1 \subset \mathcal{G}_2$ then

$$\ell(G_0|\mathcal{G}_1) \leq \ell(G_0|\mathcal{G}_2)$$

For location-scale Gaussian e-mixtures

$$\ell(G_0|\mathcal{E}_{k_0})=0.$$

For any o-mixtures: G_0 has k_0 support points, but we consider the space \mathcal{O}_k of all measures with at most $k > k_0$ support points. Then,

 $\ell(G_0|\mathcal{O}_k) \geq 1.$

In fact, for location-scale Gaussian o-mixtures, we can show

 $\ell(G_0|\mathcal{O}_k) \geq 3.$

Image: A image: A

Singularities in o-mixtures

$\mathcal{G}_0 \text{ has } k_0 \text{ support points,}$ $\mathcal{G}:=\mathcal{O}_k\text{, are space of mixing measures having at most } k>k_0 \text{ support points}$

→ ∃ →

When do we have $\ell(G_0|\mathcal{O}_k) = 1$?

Definition

[Chen, 1995; Nguyen, 2013]

G is non-singular in a o-mixture model with at most k components if

$$\left\{f(x| heta_i), rac{\partial f}{\partial heta}(x| heta_i), rac{\partial^2 f}{\partial heta^2}(x| heta_i), |i=1,\dots,k
ight\}$$

are linearly independent functions of x

Theorem

[Nguyen, 2013; Chen, 1995]

・ 何 ト ・ ヨ ト ・ ヨ ト

Under non-singularity of G_0 , and compactness conditions on parameter space, there holds

$$V(p_G,p_{G_0})\gtrsim W_2^2(G,G_0).$$

This is a corollary of our theory, due to the fact that one can establish

$$\ell(G_0|\mathcal{O}_k)=1.$$

Since MLE or Bayes estimators yield root-n density estimation rate, it implies that

$$W_2(\hat{G}_n, G_0) = O_p(n^{-1/4})$$

where G_0 denotes true mixing distribution, G_n estimate from an *n*-iid sample

イロト イポト イヨト イヨ

Since MLE or Bayes estimators yield root-*n* density estimation rate, it implies that

$$W_2(\hat{G}_n, G_0) = O_p(n^{-1/4})$$

where G_0 denotes true mixing distribution, G_n estimate from an *n*-iid sample

This result is applicable to location Gaussian o-mixture, scale Gaussian o-mixture, but not applicable to

- location-scale Gaussian o-mixture
- skewnormal o-mixture

< □ > < 同 > < 回 > < 回 > < 回 >

Location-scale Gaussian o-mixture

- given *n*-iid sample from a location-scale Gaussian mixture with mixing measure G₀ (which has k₀ components)
- fit the data with a mixture model with $k > k_0$ components
- $\ell(G_0|\mathcal{O}_k)$ is determined by $(k k_0)$ specifically, by the following system of polynomial equations:

▲御▶ ▲陸▶ ▲臣

Location-scale Gaussian o-mixture

- given *n*-iid sample from a location-scale Gaussian mixture with mixing measure G₀ (which has k₀ components)
- fit the data with a mixture model with $k > k_0$ components
- $\ell(G_0|\mathcal{O}_k)$ is determined by $(k k_0)$ specifically, by the following system of polynomial equations:

$$\sum_{j=1}^{k-k_0+1} \sum_{n_1+2n_2=\alpha} \frac{c_j^2 a_j^{n_1} b_j^{n_2}}{n_1! n_2!} = 0 \quad \text{for each } \alpha = 1, \dots, r$$
(1)

there are r equations for $3(k - k_0 + 1)$ unknowns $(c_j, a_j, b_j)_{j=1}^{k-k_0+1}$

イロト イポト イヨト イヨト 二日

Location-scale Gaussian o-mixture

- given *n*-iid sample from a location-scale Gaussian mixture with mixing measure G₀ (which has k₀ components)
- fit the data with a mixture model with $k > k_0$ components
- $\ell(G_0|\mathcal{O}_k)$ is determined by $(k k_0)$ specifically, by the following system of polynomial equations:

$$\sum_{j=1}^{k-k_0+1} \sum_{n_1+2n_2=\alpha} \frac{c_j^2 a_j^{n_1} b_j^{n_2}}{n_1! n_2!} = 0 \quad \text{for each } \alpha = 1, \dots, r$$
(1)

there are r equations for $3(k - k_0 + 1)$ unknowns $(c_j, a_j, b_j)_{j=1}^{k-k_0+1}$

 let r
 ≥ 1 minimum value of r ≥ 1 such that the above equations do not have any non-trivial real-valued solution. A solution is considered non-trivial if all c_is are non-zeros, and at least one of the a_is is non-zero.

イロト 不得 トイヨト イヨト 二日

Theorem [Singularity in localion-scale Gaussian o-mixtures] $\ell(G_0|\mathcal{O}_k) = \overline{r} - 1.$

Corrolary [Location-scale Gaussian finite mixtures]

convergence rate of mixing measure G by either MLE or Bayesian estimation is $(\log n/n)^{1/2\overline{r}}$, under both $W_{\overline{r}}$ and W_1 .

[Ho & Nguyen (2016)]

イロト イポト イヨト イヨト

More on system of r polynomial equations (1):

let us consider the case $k = k_0 + 1$, and let r = 3, then we have

$$\begin{split} c_1^2 a_1 + c_2^2 a_2 &= 0, \\ \frac{1}{2} (c_1^2 a_1^2 + c_2^2 a_2^2) + c_1^2 b_1 + c_2^2 b_2 &= 0, \\ \frac{1}{3!} (c_1^2 a_1^3 + c_2^2 a_2^3) + c_1^2 a_1 b_1 + c_2^2 a_2 b_2 &= 0 \end{split}$$

a non-trivial solution exists, by choosing $c_2=c_1\neq 0$, $b_1=b_2=-1/2$, and $a_1=1,a_2=-1$. Hence, $\overline{r}\geq 4$.

イロト イポト イヨト イヨト 二日

More on system of r polynomial equations (1):

let us consider the case $k = k_0 + 1$, and let r = 3, then we have

$$\begin{split} &c_1^2 a_1 + c_2^2 a_2 = 0, \\ &\frac{1}{2} (c_1^2 a_1^2 + c_2^2 a_2^2) + c_1^2 b_1 + c_2^2 b_2 = 0, \\ &\frac{1}{3!} (c_1^2 a_1^3 + c_2^2 a_2^3) + c_1^2 a_1 b_1 + c_2^2 a_2 b_2 = 0. \end{split}$$

a non-trivial solution exists, by choosing $c_2 = c_1 \neq 0$, $b_1 = b_2 = -1/2$, and $a_1 = 1, a_2 = -1$. Hence, $\overline{r} \geq 4$.

For r = 4, the system consists of the three equations above, plus

$$\frac{1}{4!}(c_1^2a_1^4+c_2^2a_2^4)+\frac{1}{2!}(c_1^2a_1^2b_1+c_2^2a_2^2b_2)+\frac{1}{2!}(c_1^2b_1^2+c_2^2b_2^2)=0.$$

This system has no non-trivial solution. So, for $k = k_0 + 1$, we have $\overline{r} = 4$.

イロト 不得 トイラト イラト 一日

Using the Groebner bases method, we can figure out the zeros of a system of real polynomial equations. So,

- (1) [Overfitting by one] if $k k_0 = 1$, then $\ell(G_0|\mathcal{O}_k) = 3$. So, MLE/Bayes estimation rate is $n^{-1/8}$.
- (2) [Overfitting by two] if $k k_0 = 2$, then $\ell(G_0|\mathcal{O}_k) = 5$, so the rate is $n^{-1/12}$.
- (3) [Overfitting by three] if $k k_0 = 3$, then $\ell(G_0|\mathcal{O}_k) \ge 6$. So, the rate is not better than $n^{-1/14}$.

イロト イポト イヨト イヨト 二日

Using the Groebner bases method, we can figure out the zeros of a system of real polynomial equations. So,

- (1) [Overfitting by one] if $k k_0 = 1$, then $\ell(G_0|\mathcal{O}_k) = 3$. So, MLE/Bayes estimation rate is $n^{-1/8}$.
- (2) [Overfitting by two] if $k k_0 = 2$, then $\ell(G_0|\mathcal{O}_k) = 5$, so the rate is $n^{-1/12}$.
- (3) [Overfitting by three] if $k k_0 = 3$, then $\ell(G_0|\mathcal{O}_k) \ge 6$. So, the rate is not better than $n^{-1/14}$.

Lessons for Gaussian location-scale mixtures

• do not overfit

イロト イボト イヨト イヨト

Using the Groebner bases method, we can figure out the zeros of a system of real polynomial equations. So,

- (1) [Overfitting by one] if $k k_0 = 1$, then $\ell(G_0|\mathcal{O}_k) = 3$. So, MLE/Bayes estimation rate is $n^{-1/8}$.
- (2) [Overfitting by two] if $k k_0 = 2$, then $\ell(G_0|\mathcal{O}_k) = 5$, so the rate is $n^{-1/12}$.
- (3) [Overfitting by three] if $k k_0 = 3$, then $\ell(G_0|\mathcal{O}_k) \ge 6$. So, the rate is not better than $n^{-1/14}$.

Lessons for Gaussian location-scale mixtures

- do not overfit
- if you must, be conservative in allowing extra mixing components

イロト 不得 トイラト イラト 一日



Figure : Location-scale Gaussian mixtures. From left to right: (1) Exact-fitted setting; (2) Over-fitted by one component; (3) Over-fitted by one component; (4) Over-fitted by two components.



Figure : MLE rates for location-covariance mixtures of Gaussians. L to R: (1) Exact-fitted: $W_1 \simeq n^{-1/2}$. (2) Over-fitted by one: $W_4 \simeq n^{-1/8}$. (3) Over-fitted by two: $W_6 \simeq n^{-1/12}$.

→ Ξ →

The direct link to algebraic geometry

Recall r-canonical form

$$\frac{p_{G_n}(x) - p_{G_0}(x)}{W_r^r(G_n, G_0)} = \sum_{l=1}^{L_r} \left(\frac{\xi_l(G_0, \Delta G_n)}{W_r^r(G, G_0)} \right) H_l(x) + o(1),$$

where

- coefficients $\xi_l(G_0, \Delta G_n)/W_r^r(G_n, G_0)$ are the ratio of two semi-polynomials of the parameter perturbation of G_n around G_0
- as $G_n \rightarrow G_0$, the collection of these ratios tend to a system of real polynomial equations

・ 何 ト ・ ヨ ト ・ ヨ ト

Singularities in e-mixtures of skewnormals



Figure 2: The singularity levels of $G_0 \in \mathcal{E}_{k_0}$. "C" stands for conformant and " \mathcal{NC} " stands for nonconformant. The leaf is arranged horizontally according to its level of singularity and vertically according to its influence on the corresponding settings of G_0 .

▲ 重 ▶ 重 ∽ ९ ୯ June 2016 44 / 50

< □ > < □ > < □ > < □ > < □ >



$$P_{1}(\eta) = \prod_{j=1}^{k} m_{j}.$$

$$P_{2}(\eta) = \prod_{1 \le i \ne j \le k} \left\{ (\theta_{i} - \theta_{j})^{2} + \left[\sigma_{i}^{2} (1 + m_{j}^{2}) - \sigma_{j}^{2} (1 + m_{i}^{2}) \right]^{2} \right\}.$$

Long Nguyen (UM)

Partition of subset NC, without Gaussian components



Figure 3: The level of singularity structure of G_0 relating to the nonconformant without symmetry e-mixtures setting. Here, " $\mathcal{N}C$ " stands for nonconformant. The leaf is arranged horizontally according to its level of singularity and vertically according to its influence on the corresponding setting. k^* is the maximum length of nonconformant homologous sets with C.1 singularity in G_0 respectively. Finally, the term $\mathbb{F}(G_0)$ is defined as in (27).

< □ > < 同 > < 回 > < 回 > < 回 >

Partition of subset NC, with some Gaussian components



Figure 4: The level of singularity structure of G_0 relating to the nonconformant with symmetry emixtures setting. Here, " \mathcal{NC} " stands for nonconformant. The leaf is arranged horizontally according to its level of singularity and vertically according to its influence on the corresponding setting. k^* is the maximum length of nonconformant homologous sets with C.1 singularity in G_0 respectively. The term $\overline{s}(G_0)$ is defined as in (27).

< □ > < □ > < □ > < □ > < □ > < □ >

Singularities in o-mixtures of skewnormal mixtures



Figure 6: The singularity level of $G_0 \in \mathcal{E}_k$ under the o-mixtures setting. "C" stands for conformant and "NC" stands for nonconformant. The leaf is arranged horizontally according to its level of singularity and vertically according to its influence on the corresponding setting of G_0 . The term $\mathbb{R}(G_0, k)$ is defined as in (43) and $\mathbb{R}_0(G_0, k)$ is defined as in (54). Note that, for the setting $P_2(\eta) = 0$, we assume G_0 has no generic components for the singlivity of presentation.

イロン イ理 とく ヨン イ ヨン

Summary

- Singularities are common in mixture models, including finite mixtures
- Beyond the singular points of Fisher information
- They are organized into levels, which subdivide the parameter space into multi-level partitions, each of which allow different minimax and MLE convergence rate
- Now that we know what the singularities are (mostly), how to go about improving the estimation algorithm both statistically and computationally?

For details, see

- Convergence of latent mixing measures in finite and infinite mixture models. (Ann. Statist., 2013)
- Convergence rates of parameter estimation in some weakly identifiable models (with N. Ho, Ann. Statist., 2016)
- Singularity structures and parameter estimation in finite mixture models (manuscript to be submitted)

.