2019 NCTS & Sinica Summer Course: Nonparametric Statistics and Geometric Estimation

Yen-Chi Chen

Department of Statistics University of Washington

Summer 2019

Introduction

Introduction to Geometric Estimation

- Geometric estimation studies the problem of estimating a geometric feature of a function (of interest).
- Often this function is the underlying probability density function (PDF) that generates our data.
- Other the functions of interest in statistics: the regression function, the difference between two densities/regression functions, conditional probability of an event.



The data can be viewed as

$$X_1, \cdots, X_n \sim p,$$

p is a probability density function.



The data can be viewed as

 $X_1, \cdots, X_n \sim p,$

p is a probability density function.



The data can be viewed as

 $X_1, \cdots, X_n \sim p,$

p is a probability density function.



The data can be viewed as

 $X_1, \cdots, X_n \sim p,$

p is a probability density function.



The data can be viewed as

 $X_1, \cdots, X_n \sim p,$

p is a probability density function.



The data can be viewed as

 $X_1, \cdots, X_n \sim p,$

p is a probability density function.



Nonparametric Density Estimation

Density Estimation: Introduction

• A statistical model views the data as random variables X_1, \dots, X_n from an unknown distribution function P(x) with a PDF p(x).

Density Estimation: Introduction

- A statistical model views the data as random variables X_1, \dots, X_n from an unknown distribution function P(x) with a PDF p(x).
- In most cases, we do not know the PDF p(x) but we want to reconstruct it from the data.
- ► The goal of density estimation is to estimate p(x) using X₁,..., X_n.
- In other words, the parameter of interest is the PDF p(x).

Nonparametric Approach: Introduction

- A common approach to estimate p(x) is to assume a parametric model such as a Gaussian and recover the parameters of the model by fitting to the data.
- However, this idea is often either too restrictive to capture the intricate structure of the PDF or computationally infeasible (in the case of mixture models).

Nonparametric Approach: Introduction

- A common approach to estimate p(x) is to assume a parametric model such as a Gaussian and recover the parameters of the model by fitting to the data.
- However, this idea is often either too restrictive to capture the intricate structure of the PDF or computationally infeasible (in the case of mixture models).
- An alternative approach is to estimate the PDF nonparametrically.
- Namely, we directly estimate the PDF without assuming a parametric form of the PDF.

In this lecture, we will focus on one particular nonparametric estimator-the kernel density estimator (KDE).

- In this lecture, we will focus on one particular nonparametric estimator-the kernel density estimator (KDE).
- ► The KDE estimates the PDF using the following form:

$$\widehat{p}_h(x) = \frac{1}{nh} \sum_{i=1}^n K\left(\frac{x-X_i}{h}\right),$$

where K(x) is a function called the kernel function and h > 0 is a quantity called smoothing bandwidth that controls the amount of smoothing.

- In this lecture, we will focus on one particular nonparametric estimator-the kernel density estimator (KDE).
- ► The KDE estimates the PDF using the following form:

$$\widehat{p}_h(x) = \frac{1}{nh} \sum_{i=1}^n K\left(\frac{x-X_i}{h}\right),$$

where K(x) is a function called the kernel function and h > 0 is a quantity called smoothing bandwidth that controls the amount of smoothing.

• Common choice of K(x) includes the Gaussian $K(x) = \frac{1}{\sqrt{2\pi}}e^{-\frac{x^2}{2}}$ and the uniform $K(x) = \frac{1}{2}I(-1 \le x \le 1)$.

- In this lecture, we will focus on one particular nonparametric estimator-the kernel density estimator (KDE).
- ► The KDE estimates the PDF using the following form:

$$\widehat{p}_h(x) = \frac{1}{nh} \sum_{i=1}^n K\left(\frac{x-X_i}{h}\right),$$

where K(x) is a function called the kernel function and h > 0 is a quantity called smoothing bandwidth that controls the amount of smoothing.

- Common choice of K(x) includes the Gaussian $K(x) = \frac{1}{\sqrt{2\pi}}e^{-\frac{x^2}{2}}$ and the uniform $K(x) = \frac{1}{2}I(-1 \le x \le 1)$.
- The idea of KDE is: we smooth out each data point using the kernel function into small bumps and then we sum over all bumps to obtain a density estimate.



Black dots: locations of observations. Purple bumps: the kernel function at each observation. Brown curve: final density estimate from KDE.



The kernel function generally does not affect the density estimate too much. $$10\,/\,67$$



The smoothing bandwidth often has a much stronger effect on the quality of estimation.

In statistics, there are two common quantities to measure the accuracy of estimation- bias and variance of an estimator.

- In statistics, there are two common quantities to measure the accuracy of estimation- bias and variance of an estimator.
- ▶ When the smoothing bandwidth $h \approx 0$ and sample size n is large,

$$\mathbf{bias}(\widehat{p}_h(x)) = \mathbb{E}(\widehat{p}_h(x)) - p(x) = C_{1,K}p''(x)h^2 + o(h^2),$$

and the variance has the asymptotic form:

$$\operatorname{Var}(\widehat{p}_h(x)) = C_{2,\kappa} \frac{p(x)}{nh^d} + o\left(\frac{1}{nh^d}\right),$$

where $C_{1,\mathcal{K}}$ and $C_{2,\mathcal{K}}$ are constants depending on the kernel function.

- In statistics, there are two common quantities to measure the accuracy of estimation- bias and variance of an estimator.
- ▶ When the smoothing bandwidth $h \approx 0$ and sample size n is large,

$$\mathbf{bias}(\widehat{p}_h(x)) = \mathbb{E}(\widehat{p}_h(x)) - p(x) = C_{1,K}p''(x)h^2 + o(h^2),$$

and the variance has the asymptotic form:

$$\operatorname{Var}(\widehat{p}_h(x)) = C_{2,\kappa} \frac{p(x)}{nh^d} + o\left(\frac{1}{nh^d}\right)$$

where $C_{1,\mathcal{K}}$ and $C_{2,\mathcal{K}}$ are constants depending on the kernel function.

 The mean squared error (MSE) is a common quantity of measuring the accuracy that takes both the bias and variance into consideration. For the KDE, it is

$$\mathsf{MSE}(\widehat{p}_h(x)) = \mathbf{bias}^2(\widehat{p}_h(x)) + \mathsf{Var}(\widehat{p}_h(x))$$
$$= C_{1,\mathcal{K}}^2 |p''(x)|^2 h^4 + C_{2,\mathcal{K}} \frac{p(x)}{nh^d} + o(h^4) + o\left(\frac{1}{nh^d}\right)/67$$

The MSE measures the accuracy at a single point x. The overall performance is often quantified by the mean integrated squared error (MISE):

$$\mathsf{MISE}(\widehat{p}_h) = \int \mathsf{MSE}(\widehat{p}_h(x))$$
$$= C_{1,K}^2 \left(\int |p''(x)|^2 dx \right) h^4 + \frac{C_{2,K}}{nh^d} + o(h^4) + o\left(\frac{1}{nh^d}\right)$$

This implies several interesting facts:

- 1. when h is too large, we suffer from the bias.
- 2. when h is too small, we suffer from the variance.
- 3. the optimal choice is $h \simeq n^{-1/(d+4)}$.

L_∞ Analysis - 1

- The MISE is essentially just the L_2 distance between \hat{p}_h and p.
- ► We can then generalize the result to other L_p distance between the two quantities.
- ► Among all p, one particularly interesting case is L_∞ distance. In this case,

$$\sup_{x} |\widehat{p}_h(x) - p(x)| = \|\widehat{p}_h - p\|_{\infty} = O(h^2) + O_P\left(\sqrt{\frac{\log n}{nh^d}}\right).$$

The above bound follows from the following decomposition:

$$\|\widehat{p}_h - p\|_{\infty} \leq \underbrace{\|\widehat{p}_h - p_h\|_{\infty}}_{O_P} + \underbrace{\|p_h - p\|_{\infty}}_{O},$$

where $p_h = \mathbb{E}(\widehat{p}_h) = p \otimes K$ is also called the smoothed density.

L_{∞} Analysis - 2

$$\|\widehat{p}_h - p\|_{\infty} \leq \underbrace{\|\widehat{p}_h - p_h\|_{\infty}}_{O_P} + \underbrace{\|p_h - p\|_{\infty}}_{O},$$

- ► The fact that the bias term at rate O(h²) is from the usual analysis (Taylor expansion).
- The bound on the stochastic variation is more involved.
- ► In short, it follows from the Talagrand's inequality:

$$P(\|\widehat{p}_h - p_h\|_{\infty} > t) \le c_0 e^{-c_1 n h^d t^2}$$

when $t > \sqrt{\frac{|\log h|}{nh^d}}$. This result is formally given in Giné and Guillou (2002).

In fact, after rescaling, the random variable ||p̂_h − p_h||_∞ converges in distribution to an extreme value distribution (Bickel and Rosenblatt (1973)).

L_{∞} Analysis - 3

- ► The L_∞ analysis can be generalized to derivatives of the density function.
- Gradient:

$$\sup_{x} \|\nabla \widehat{p}_h(x) - \nabla p(x)\|_{\max} = O(h^2) + O_P\left(\sqrt{\frac{\log n}{nh^{d+2}}}\right).$$

Hessian:

$$\sup_{x} \|\nabla \nabla \widehat{p}_{h}(x) - \nabla \nabla p(x)\|_{\max} = O(h^{2}) + O_{P}\left(\sqrt{\frac{\log n}{nh^{d+4}}}\right).$$

And higher-order derivatives can be derived accordingly.

Application of L_∞ Analysis - 1

- ► The L_∞ analysis implies a construction of a simultaneous confidence band of p.
- There are two types of confidence bands for a function p: pointwise confidence bands and simultaneous confidence bands.

Application of L_∞ Analysis - 1

- ► The L_∞ analysis implies a construction of a simultaneous confidence band of p.
- There are two types of confidence bands for a function p: pointwise confidence bands and simultaneous confidence bands.
- Pointwise CI: for any given point x and confidence level 1 − α, we construct an interval C_{1−α} = [ℓ_{1−α}, u_{1−α}] from the data such that

$$P(\ell_{1-\alpha} \leq p(x) \leq u_{1-\alpha}) = 1 - \alpha + o(1).$$

Simultaneous CB (confidence band): given 1 − α, we construct a band C_{1−α}(x) = [L_{1−α}(x), U_{1−α}(x)] from the data such that

$$P(L_{1-\alpha}(x) \leq p(x) \leq U_{1-\alpha}(x) \text{ for all } x) = 1 - \alpha + o(1).$$

Application of L_{∞} Analysis - 2

We can construct a simultaneous confidence band by bootstrapping the L_{∞} distance.



Pointwise CI (left) and simultaneous CB (right)¹.

¹More details can be found in: https://arxiv.org/abs/1702.07027

Application of L_{∞} Analysis - 3

- ► The L_∞ analysis also implies the convergence of geometric structures.
- In particular, some geometric structures converge when the derivatives converge².
- Convergence rate depending on $\|\widehat{p}_h p_h\|_{\infty}$:
 - Level sets, cluster trees, and persistent diagrams.
- Convergence rate depending on $\|\nabla \hat{p}_h \nabla p_h\|_{\infty}$:
 - ► Local modes, Morse-Smale complex, and gradient system.
- Convergence rate depending on $\|\nabla \nabla \widehat{p}_h \nabla \nabla p_h\|_{\infty}$:

Ridges.

²A tutorial on this topic is in: https://arxiv.org/abs/1704.03924

Geometric Estimation

Level Sets - 1

• Given a level $\lambda > 0$, the density level set is

$$L_{\lambda} = \{ x : p(x) = \lambda \}.$$

• A natural estimator of L_{λ} is the plug-in using a KDE³:

$$\widehat{L}_{\lambda} = \{ x : \widehat{p}_h(x) = \lambda \}.$$

Note that sometime in the literature, the set of interest is the upper level set:

$$S_{\lambda} = \{x : p(x) \geq \lambda\}.$$

Under smoothness conditions, the boundary of the upper level set is the level set L_{λ} .

 Level set is a particularly interesting example so we take a deeper look at this problem.

³Materials on this topic can be found in:

https://arxiv.org/abs/1504.05438 and the reference therein.

Level Sets - 2

- Level set has two common applications:
 - 1. Anomaly detection-observations in the low density area (outside the level set) are going to be classified as anomaly.
 - 2. Clustering-observations inside the level set (high density area) are going to be clustered together.


- There has been a tremendous amount of literature on the convergence of level set.
- Often the convergence is expressed in terms of the Hausdorff distance

$$Haus(A, B) = \max\{\sup_{x \in A} d(x, B), \sup_{x \in B} d(x, A)\},\$$

where $d(x, A) = \inf_{y \in A} ||x - y||$ is the distance from a point x to a set A.

► The Hausdorff distance can be viewed as an L_∞ distance for sets.

A common assumption to ensure the convergence of Hausdorff distance is the gradient bound:

$$\inf_{x\in L_{\lambda}} \|\nabla p(x)\| \geq g_0 > 0,$$

for some constant g_0 .

 Under this assumption (and some other common assumptions),

$$\mathsf{Haus}(\widehat{\mathcal{L}}_{\lambda},\mathcal{L}_{\lambda})=O_{\mathsf{P}}\left(\|\widehat{p}_{\mathsf{h}}-p\|_{\infty}
ight).$$

- If we further assume that p(x) has bounded second derivative everywhere, then the level set L_λ is smooth in the sense that the *reach* is positive.
- The reach of a set is the longest distance away from a set that still has a unique projection back to the set.
- If L_λ has a reach r₀, then for any point x with d(x, L_λ) < r₀, x has a unique projection back to L_λ.
- ► The positive reach properties of level sets imply that \widehat{L}_{λ} and L_{λ} are (asymptotically) *normal compatible*, meaning that there is a unique projection from every point in L_{λ} to \widehat{L}_{λ} and vice versa.

 The normal compatibility implies that we can decompose the Hausdorff distance as

$$\mathsf{Haus}(\widehat{L}_{\lambda}, L_{\lambda}) = \sup_{x \in L_{\lambda}} d(x, \widehat{L}_{\lambda}).$$

• An more interesting fact is that for any $x \in L_{\lambda}$,

$$d(x, \widehat{L}_{\lambda}) = rac{1}{\|
abla p(x)\|} |\widehat{p}_h(x) - p(x)| + ext{smaller order terms}.$$

This implies that asymptotically,

$$\mathsf{Haus}(\widehat{L}_{\lambda}, L_{\lambda}) = \sup_{x \in L_{\lambda}} \frac{1}{\|\nabla p(x)\|} |\widehat{p}_{h}(x) - p(x)|,$$

which is the supremum of a stochastic process defined over the manifold.

$$\mathsf{Haus}(\widehat{L}_{\lambda}, L_{\lambda}) = \sup_{x \in L_{\lambda}} \frac{1}{\|\nabla p(x)\|} |\widehat{p}_{h}(x) - p(x)|.$$

- This shows that the Hausdorff distance follows an extreme value distribution after rescaling.
- Also, it implies that we can use the bootstrap to construct a confidence set of L_λ.



Another interesting geometric structure is the local modes of the PDF:

$$M = \{x : \nabla p(x) = 0, \lambda_1(x) < 0\},\$$

where $\lambda_1(x)$ is the largest eigenvalue of the Hessian matrix $\nabla \nabla p(x)^4$.

Similar to the level set problem, a simple estimator of M is the plug-in from the KDE

$$\widehat{M}_h = \{ x : \nabla \widehat{p}_h(x) = 0, \widehat{\lambda}_1(x) < 0 \}.$$

⁴A tutorial on this topic is in: https://arxiv.org/abs/1406.1780 and https://arxiv.org/abs/1408.1381.

- Before talking about the applications of local modes, we first discuss some of its properties.
- If the density function p is a Morse function, i.e., all critical points of p are non-degenerated, then
 - 1. $\operatorname{Haus}(\widehat{M}_h, M) = O_P\left(\sup_x \|\nabla \widehat{p}_h(x) \nabla p(x)\|_{\max}\right),$
 - 2. with a probability approaching to 1, there exists a one-to-one correspondence between elements in \widehat{M}_h and elements in M.
- ► A common assumption to replace the Morse condition is that there exists a lower bound <u>λ</u> > 0 such that

$$\min_{x\in M} |\lambda_1(x)| \geq \underline{\lambda} > 0.$$

In fact, one can obtain a faster convergence rate of Haus(M̂_h, M) without the log n term in the variance.

- Local modes can be used to perform a cluster analysis.
- This is known as the mode clustering method (mean-shift clustering).



- ► The clusters are defined through a gradient system of p (or in the sample case, p̂_h).
- For a given point x, we define a gradient flow π_x such that

$$\pi_x(0) = x, \quad \pi'_x(t) = \nabla p(\pi_x(t)).$$

- The destination π_x(∞) = lim_{t→∞} π_x(t) ∈ M when p is a Morse function for almost every x except a set of point with Lebesgue measure 0.
- ► Thus, we can use the destination π_x(∞) of each x to cluster data points. Namely, points with the same destination will be clustered together.



- Numerically, we use the *mean-shift* algorithm to do the gradient ascent.
- Let x_0 be the initial point.
- We iterate

$$x_{t+1} = \frac{\sum_{i=1}^{n} X_i K\left(\frac{X_i - x_t}{h}\right)}{\sum_{i=1}^{n} K\left(\frac{X_i - x_t}{h}\right)}$$

until convergence.

 Note that this works for Gaussian kernel. Some other kernel functions also work after modifications.

- For each m ∈ M, let D(m) = {x : π_x(∞) = m} be the basin of attraction with respect to m.
- The set D = {D(m) : m ∈ M} forms a partition of the entire support of p (except for a set of Lebesgue measure 0).
- ▶ Similarly, we may define the sample version of it $\{\widehat{D}(m) : m \in \widehat{M}_h\}$, where $\widehat{D} = \widehat{D}(m) = \{x : \widehat{\pi}_x(\infty) = m\}$ is the basin of attraction using the gradient system of \widehat{p}_h .

⁵Materials on this topic can be found in: https://arxiv.org/abs/1506.08826

- Let B = {∂D(m) : m ∈ M} be the collection of boundaries of the basins and B = {∂D(m) : m ∈ M} be the sample version of it.
- ► The convergence of D toward D can be characterized by the convergence of B to B.
- If for every x ∈ M, p(x) are convex with respect to the 'normal space' of B at x, then

$$\operatorname{Haus}(\widehat{\mathcal{B}},\mathcal{B}) = O_P\left(\sup_{x} \|\nabla \widehat{p}_h(x) - \nabla p(x)\|_{\max}\right).$$

- There is an elegant idea combining both level sets and modes: cluster tree⁶.
- The cluster tree considers the collection of clusters formed by the upper level sets and keeps track of their relationships.
- When applying to a density function, a cluster tree is also called a density tree.

⁶Materials on this topic can be found in: https://arxiv.org/abs/1605.06416











































- Level sets are the basis of constructing a cluster tree.
- Local modes are associated to the creation of a new branch in a cluster tree.
- Saddle points or local minima are related to the elimination (merging) of a branch in a cluster tree.
- Cluster tree provides an elegant way to represent the shape of the PDF and can be used to visualize the data.
Let T_p be the cluster tree based on the PDF and $\hat{T}_p = T_{\hat{p}_h}$ be the cluster tree based on the KDE.

To measure the estimation error, a simple metric is

$$d_{\infty}(\widehat{T_p}, T_p) = \sup_{x} \|\widehat{p}_h(x) - p(x)\|,$$

so the convergence rate is

$$d_{\infty}(\widehat{T_p},T_p)=O(h^2)+O_P\left(\sqrt{rac{\log n}{nh^d}}
ight).$$

Let T_p be the cluster tree based on the PDF and $\hat{T}_p = T_{\hat{p}_h}$ be the cluster tree based on the KDE.

To measure the estimation error, a simple metric is

$$d_{\infty}(\widehat{T_p}, T_p) = \sup_{x} \|\widehat{p}_h(x) - p(x)\|,$$

so the convergence rate is

$$d_{\infty}(\widehat{T_p},T_p)=O(h^2)+O_P\left(\sqrt{rac{\log n}{nh^d}}
ight).$$

 Another way of defining statistical convergence is based on the probability

$$P_n = P\left(\widehat{T_p} \text{ and } T_p \text{ are topological equivalent}\right).$$

Let T_p be the cluster tree based on the PDF and $\hat{T}_p = T_{\hat{p}_h}$ be the cluster tree based on the KDE.

To measure the estimation error, a simple metric is

$$d_{\infty}(\widehat{T_p}, T_p) = \sup_{x} \|\widehat{p}_h(x) - p(x)\|,$$

so the convergence rate is

$$d_\infty(\widehat{T_p},T_p)=O(h^2)+O_P\left(\sqrt{rac{\log n}{nh^d}}
ight).$$

 Another way of defining statistical convergence is based on the probability

$$P_n = P\left(\widehat{T_p} \text{ and } T_p \text{ are topological equivalent}\right).$$

• Under smoothness conditions and $n \to \infty, h \to 0$,

$$P_n \geq 1 - e^{-nh^{d+4} \cdot C_p},$$

for some constant C_p depending on the density function p.

 There are other notions of convergence/consistency of a tree estimator.

- There are other notions of convergence/consistency of a tree estimator.
- Convergence in the merge distortion metric (Eldridge et al. 2015) is one example.

- There are other notions of convergence/consistency of a tree estimator.
- Convergence in the merge distortion metric (Eldridge et al. 2015) is one example.
- ▶ However, it was shown in Kim et al. (2016) that this metric is equivalent to the L_{∞} metric.

- There are other notions of convergence/consistency of a tree estimator.
- Convergence in the merge distortion metric (Eldridge et al. 2015) is one example.
- ▶ However, it was shown in Kim et al. (2016) that this metric is equivalent to the L_{∞} metric.
- Hartigan consistency (Chaudhuri and Dasgupta 2010; Balakrishnan et al. 2013) is another way to measure the consistency of a tree estimator.

- There are other notions of convergence/consistency of a tree estimator.
- Convergence in the merge distortion metric (Eldridge et al. 2015) is one example.
- ▶ However, it was shown in Kim et al. (2016) that this metric is equivalent to the L_{∞} metric.
- Hartigan consistency (Chaudhuri and Dasgupta 2010; Balakrishnan et al. 2013) is another way to measure the consistency of a tree estimator.
- Note: cluster tree can also be recovered by a kNN approach; see Chaudhuri and Dasgupta (2010) and Chaudhuri et al. (2014) for more details.

Persistent Diagrams - 1

- Cluster trees contain only the information about connected components of level sets.
- Connected components are 0th order homology group.
- One can generalize this concept to higher order homology groups.
- The creation and elimination of homology groups can be summarized using the persistent diagram.

Persistent Diagrams - 2

- Again, we may define the persistent diagram formed by using the population PDF p and the KDE p̂_h.
- ▶ Let PD = PD(p) be the persistent diagram formed by level sets of p and $\widehat{PD} = PD(\widehat{p}_h)$ be the one formed by level sets of \widehat{p}_h .
- Using the fact the bottleneck distance of persistent diagrams is bounded by the L_{∞} distance of the generated function⁷, we conclude that

$$d_B(\widehat{\mathsf{PD}},\mathsf{PD}) \le \|\widehat{p}_h - p\|_{\infty} = O(h^2) + O_P\left(\sqrt{rac{\log n}{nh^d}}\right)$$

⁷https://link.springer.com/article/10.1007/s00454-006-1276-5.

Persistent Diagrams - 3

- There are other distance for persistent diagrams such as the Wasserstein distance.
- ▶ But the bottleneck distance has a nice property that it is an L_∞ type distance so we can use it to construct a confidence set⁸.
- ► This is often done by bootstrapping the upper bound ||p̂_h - p||_∞ since it is unclear if bootstrapping the bottleneck distance will work or not.
- ► Also, computing the bottleneck distance is challenging.

⁸See https://arxiv.org/abs/1303.7117 for more details.

- Ridges are another interesting geometric structure that we may want to study⁹.
- They can be viewed as generalized local modes.



⁹Materials on this topic can be found in: https://arxiv.org/abs/1406.5663 and https://arxiv.org/abs/1212.5156

- Here is the formal definition of ridges.
- Let v₁(x), · · · , v_d(x) be the ordered eigenvectors of ∇∇p(x), where v₁(x) corresponds to the largest eigenvalue.
- Define $V(x) = [v_2(x), \cdots, v_d(x)] \in \mathbb{R}^{d \times (d-1)}$.
- Ridges are defined as the collection:

$$R = \{x : V(x)V(x)^{T}\nabla p(x) = 0, \lambda_{2}(x) < 0\}.$$

- V(x)V(x)^T is the projection matrix onto the subspace spanned by v₂(x), · · · , v_d(x).
- ▶ Thus, the ridge *R* is the collection of *projected local modes*.



An application of ridges in Astronomy.



An application of ridges in Astronomy.

- Ridges can be estimated by the KDE.
- Let V
 _h(x) be the KDE version of V(x) and λ
 ₂(x) be the KDE version of λ₂(x). Then the ridge estimator is

$$\widehat{R}_h = \{ x : \widehat{V}_h(x) \widehat{V}_h(x)^T \nabla \widehat{p}_h(x) = 0, \widehat{\lambda}_2(x) < 0 \}.$$

- One can use the subspace constrained mean shift algorithm¹⁰ to numerically calculate the estimator.
- The convergence rate is

$$\mathsf{Haus}(\widehat{R}_h, R) = O(\sup_{x} \|\nabla \nabla \widehat{p}_h(x) - \nabla \nabla p(x)\|_{\mathsf{max}}).$$

¹⁰See http://www.jmlr.org/papers/v12/ozertem11a.html.

Singular distribution

Failure of the KDE - 1

- In the previous few sections, we see that the KDE is a powerful tool.
- However, it may not work in certain situations.



Failure of the KDE - 2

- The KDE does not work because there is no underlying PDF for GPS data!
- A better model to describe a GPS data is the following distribution:

$$P_{\text{GPS}}(x) = \pi_0 P_0(x) + \pi_1 P_1(x) + \pi_2 P_2(x),$$

where $P_0(x)$ is a distribution of point mass, and $P_1(x)$ is a distribution of a 1D density function, and $P_2(x)$ is a distribution of a 2D density function, and $\pi_0 + \pi_1 + \pi_2 = 1$ with $\pi_j \ge 0$ are proportions.

The components P₀ and P₁ make the distribution function singular so the KDE diverges.

- Although the KDE fails, the *density ranking* still works¹¹!
- The density ranking is a density surrogate that

$$\widehat{\alpha}_h(x) = \frac{1}{n} \sum_{i=1}^n I(\widehat{p}_h(x) \ge \widehat{p}_h(X_i)).$$

- ► It preserves the ordering of p̂_h(x) on each observed data points.

•
$$\widehat{\alpha}_h(x) \in [0,1]$$
 so it will not diverge.

¹¹See https://arxiv.org/abs/1611.02762 and https://arxiv.org/abs/1708.05017.

Another example of possibly singular distribution in an Astronomy dataset.



- To see why the density ranking is stable in handling GPS data (or more general, singular distributions), we consider its population quantity.
- When the PDF exists, it is easy to see that the population quantity is

$$\alpha(x) = P(p(x) \ge p(X)),$$

where $X \sim P$.

When the distribution is singular, we need to use the concept of *Hausdorff density*.

- Let C_d be the volume of a d dimensional unit ball and $B(x,r) = \{y : ||x y|| \le r\}.$
- For any integer *s*, we define

$$\mathcal{H}_{s}(x) = \lim_{r \to 0} \frac{P(B(x,r))}{C_{s}r^{s}}.$$

- → H_s(x) occurs in three regimes: 0, ∞, or a number between (0,∞).
- ► Example of 0: s = 1 on a place with 2D density (s < the structural dimension).</p>
- ► Example of ∞: s = 1 on a point mass (s > the structural dimension).
- For a point x, we then define

$$au(x) = \max\{s \leq d : \mathcal{H}_s(x) < \infty\}, \quad \rho(x) = \mathcal{H}_{\tau(x)}(x).$$

Density Ranking: Example - 1

► Assume the distribution function P is a mixture of a 2D uniform distribution within [-1, 1]², a 1D uniform distribution over the ring {(x, y) : x² + y² = 0.5²}, and a point mass at (0.5, 0), then the support can be partitioned as follows:



Density Ranking: Example - 2



- Orange region: $\tau(x) = 2 \Leftrightarrow$ contribution of $P_2(x)$.
- Red region: $\tau(x) = 1 \Leftrightarrow$ contribution of $P_1(x)$.
- Blue region: $\tau(x) = 0 \Leftrightarrow$ contribution of $P_0(x)$.

Hausdorff Density and Ranking - 1

- The function $\tau(x)$ measures the dimension of P at point x.
- ► The function p(x) describes the density of that corresponding dimension.
- We can use τ and ρ to compare any pairs of points and construct a ranking.
- ▶ For two points x_1, x_2 , we define an ordering such that $x_1 \succ_{\tau,\rho} x_2$ if

$$au(x_1) < au(x_2), \qquad ext{or} \quad au(x_1) = au(x_2), \quad
ho(x_1) >
ho(x_2).$$

Namely, we first compare the dimension of the two points, the lower dimensional structure wins. If they are on regions of the same dimension, we then compare the density of that dimension. Hausdorff Density and Ranking - 2

► Using the ordering ≻_{τ,ρ}, we then define the population density ranking as

$$\alpha(x) = P(x \succeq_{\tau,\rho} X_1)$$

• When the PDF exists, the ordering $\succ_{\tau,\rho}$ equals to $\succ_{d,p}$ so

$$\alpha(x) = P(x \succeq_{d,p} X_1) = P(p(x) \ge p(X_1)),$$

which recovers the definition in the simple case.

- In singular measure, there is a new type of critical points. We call them the *dimensional critical points*.
- These critical points contribute to the change of topology of level sets as the usual critical points but they cannot be defined by setting gradient to be 0.

- The box in the following figure is a dimensional critical point.
- Note: this is a mixture of 2D distribution and a 1D distribution on the black line (maximum value occurs at the cross).



- The box in the following figure is a dimensional critical point.
- Note: this is a mixture of 2D distribution and a 1D distribution on the black line (maximum value occurs at the cross).



- The box in the following figure is a dimensional critical point.
- Note: this is a mixture of 2D distribution and a 1D distribution on the black line (maximum value occurs at the cross).



- The box in the following figure is a dimensional critical point.
- Note: this is a mixture of 2D distribution and a 1D distribution on the black line (maximum value occurs at the cross).



Convergence under Singular Measure: Density Ranking - 1

 When P is a singular distribution and satisfies certain regularity conditions,

$$\int |\widehat{\alpha}_h(x) - \alpha(x)|^2 \, dP(x) = O(h) + O_P\left(\sqrt{\frac{1}{nh^d}}\right)$$

• Intuition of convergence: as $h \rightarrow 0$, the KDE

$$\widehat{p}_h(x) = \frac{1}{nh^d} \sum_{i=1}^n K\left(\frac{X_i - x}{h}\right)$$

diverges when x is in a lower dimensional structure $(\tau(x) < d)$.

- The bias of order O(h) is due to the smoothing from a nearby lower dimensional structure.
- However, the speed of diverging depends on $\tau(x)$. The smaller $\tau(x)$, the faster (actually the diverging rate is $O(h^{\tau(x)-d})$).
- Note: we do not estimate $\tau(x)$ when using $\widehat{\alpha}_h(x)$!

Convergence under Singular Measure: Density Ranking - 2

- Although we have L₂(P) convergence (also we have L₂ and pointwise convergence), we do not have a uniform convergence.
- Example of non-convergence of supreme norm: consider a sequence of points on a higher dimensional space but moving toward a lower dimensional structure within distance ^h/₂.
- Interestingly, we can still prove that some topological features (local modes, level sets, cluster trees, persistent diagrams) are converging.

Convergence under Singular Measure: Density Ranking - 3

- Although we do not have uniform convergence, many geometric structures still converge.
- The cluster tree, local modes, and 0-th order homology groups of density ranking converge to the population version.
- However, it is unclear if other quantities such as ridges or higher-order homology groups converge.

Conclusion

- In this lecture, we study the problem of estimating a geometric structure of the underlying PDF.
- We show that we can estimate the PDF, level sets, local modes, gradient systems, cluster trees, persistent diagrams, and ridges using the KDE.
- There are two alternative ways to view this problem.
- First, when the KDE converges to the true population PDF, many geometric structures also converge.
- Second, many geometric structures have an intrinsic stability with respect to the underlying function so small perturbation will not change it drastically.
Useful References

- Wasserman, Larry. "Topological data analysis." Annual Review of Statistics and Its Application 5 (2018): 501-532.
- Chen, Yen-Chi. "A tutorial on kernel density estimation and recent advances." Biostatistics & Epidemiology 1, no. 1 (2017): 161-187.
- ► All of statistics: a concise course in statistical inference. Larry Wasserman. Springer Science & Business Media, 2013.
- ► All of nonparametric statistics. Larry Wasserman. Springer, 2006.
- Multivariate density estimation: theory, practice, and visualization. David Scott. John Wiley & Sons, 2015.