The Intersection of Statistics and Topology: Confidence Sets for Persistence Diagrams

Brittany Terese Fasy

joint work with S. Balakrishnan, F. Lecci, A. Rinaldo, A. Singh, L. Wasserman

3 June 2014

[**F**LRWBS] Statistical Inference for Persistent Homology: Confidence Sets for Persistence Diagrams. ArXiv 1303.7117. Tentatively accepted, Annals of Statistics.

B. Fasy (Tulane)

Statistics and Topology

The Intersection of Statistics and Topology: Stochastic Convergence of Persistence Landscapes and Silhouettes

Brittany Terese Fasy

joint work with F. Chazal, F. Lecci, B. Michel A. Rinaldo, L. Wasserman

3 June 2014

[CFLRW] Stochastic Convergence of Persistence Landscapes and Silhouettes. ArXiv. SoCG 2014 Proceedings, Kyoto. Journal version pending. The Intersection of Statistics and Topology: Confidence Sets for Persistence Diagrams

Brittany Terese Fasy

joint work with S. Balakrishnan, F. Lecci, A. Rinaldo, A. Singh, L. Wasserman

3 June 2014

[FLRWBS] Statistical Inference for Persistent Homology: Confidence Sets for Persistence Diagrams. ArXiv 1303.7117. Tentatively accepted, Annals of Statistics.

B. Fasy (Tulane)

Statistics and Topology



















An Informal Discussion

A *homeomorphism* is a continuous bending, stretching, or shrinking (but not tearing or glueing) of one object into another object.

An Informal Discussion

A *homeomorphism* is a continuous bending, stretching, or shrinking (but not tearing or glueing) of one object into another object.

Circle = Square



An Informal Discussion

A *homeomorphism* is a continuous bending, stretching, or shrinking (but not tearing or glueing) of one object into another object.

Circle = Square





An Informal Discussion

A *homeomorphism* is a continuous bending, stretching, or shrinking (but not tearing or glueing) of one object into another object.

Circle = Square







Nerve Complexes



Nerve Complexes



Nerve Complexes



Nerve Complexes



Nerve Lemma

If X is a set of closed convex sets, then Nrv(X) is topologically equivalent to $\bigcup X$ (up to homotopy type).



Simplicial Homology

• X_p is the power set of all *p*-simplices.

•
$$C_p = (X_p, +_2).$$

• $\partial_{\rho}: C_{\rho} \to C_{\rho-1}$ is the boundary map.

•
$$H_p = \operatorname{Ker}(\partial_p) / \operatorname{Im}(\partial_{p+1})$$

•
$$\beta_p = \operatorname{rank}(H_p)$$
.



$$\beta_0 = n \qquad \qquad \beta_0 = 5$$
$$\beta_1 = 0 \qquad \qquad \beta_1 = 1$$

$$\beta_1 = 0$$
 $\beta_1 =$







Lower-Level Set Filtration



Lower-Level Set Filtration



Lower-Level Set Filtration



Lower-Level Set Filtration



Lower-Level Set Filtration



Lower-Level Set Filtration



Lower-Level Set Filtration



Lower-Level Set Filtration



Lower-Level Set Filtration



Upper-Level Set Filtration



Distance and Stability

A Metric on Persistence Diagrams


Distance and Stability

A Metric on Persistence Diagrams



Bottleneck Distance

Given two persistence diagrams, find the best *perfect matching M* between the point sets.

Minimize Cost

We wish to find

$$W_{\infty} = \min_{M} \{ \max_{(p,q) \in M} ||p-q||_{\infty} \}.$$



Bottleneck Distance



Minimize Cost

We wish to find

$$W_{\infty} = \min_{M} \{ \max_{(p,q) \in M} ||p-q||_{\infty} \}.$$



Stability of Matchings



Topological Inference



B. Fasy (Tulane)

Objective

Let $\mathcal{D}_{\mathcal{T}}$ denote the set of all \mathcal{T} -bounded persistence diagrams.

Confidence Sets

Given $\alpha \in (0,1)$ and unknown diagram D, we want $\mathcal{C}_{\alpha} \subset \mathcal{D}_{\mathcal{T}}$ such that

 $\mathbb{P}(D \in \mathcal{C}_{\alpha}) \geq 1 - \alpha.$

Question

If \widehat{D} is an estimate of D, how close is \widehat{D} to D?

Confidence Sets for Persistent Diagrams



Confidence Sets for Persistent Diagrams





Computing a Confidence Interval With Infinite Resources

Repeatedly sample *n* data points, obtaining:



Computing a Confidence Interval With Infinite Resources

Repeatedly sample *n* data points, obtaining: $\hat{\Theta}_{n,1}, \ldots, \hat{\Theta}_{n,N}$



Computing a Confidence Interval With Infinite Resources

Repeatedly sample *n* data points, obtaining: $\hat{\Theta}_{n,1}, \ldots, \hat{\Theta}_{n,N}$ via simulation.



When We Can Only Take One Sample

We have one sample: $S_n = \{X_1, \ldots, X_n\}$

When We Can Only Take One Sample

We have one sample: $S_n = \{X_1, \ldots, X_n\}$

Subsample (with replacement), obtaining: $\{X_1^*, \ldots, X_n^*\}$

When We Can Only Take One Sample

We have one sample: $S_n = \{X_1, \ldots, X_n\}$

Subsample (with replacement), obtaining: $\{X_1^*, \ldots, X_n^*\}$

Compute
$$\hat{\Theta}_n^* = \Theta(X_1^*, \dots, X_n^*).$$

When We Can Only Take One Sample

We have one sample: $S_n = \{X_1, \ldots, X_n\}$

Subsample (with replacement), obtaining: $\{X_1^*, \ldots, X_n^*\}$

Compute
$$\hat{\Theta}_n^* = \Theta(X_1^*, \dots, X_n^*).$$

Repeat *N* times, obtaining: $\hat{\Theta}_{n,1}^*, \dots, \hat{\Theta}_{n,N}^*$.

When We Can Only Take One Sample



Distance to a Subset



Distance to a Subset

$$egin{aligned} &d_{\mathbb{M}}(a) = \inf_{x \in \mathbb{M}} \left| \left| x - a
ight| \ &D = \mathsf{Dgm}_p^-(d_{\mathbb{M}}) \end{aligned}$$

P has continuous density p.support(P) = M. $S_n = \{X_1, \dots, X_n\} \sim P$ $\widehat{D} = \mathsf{Dgm}_p^-(d_{S_n})$



Subsampling

A Variant of the Bootstrap

Let S_b^i be a subsample of size b < n, for i = 1, ..., B. L_b is the CDF of $H(S_b^i, S_n)$.

Subsampling

A Variant of the Bootstrap

Let S_b^i be a subsample of size b < n, for i = 1, ..., B. L_b is the CDF of $H(S_b^i, S_n)$.

Confidence Sets from Subsampling [FLRWBS]

Assume that p(x) is bounded away from zero. Then, almost surely, for all large n,

$$\mathbb{P}\left(||\textit{d}_{\mathbb{M}}-\textit{d}_{\mathcal{S}_n}||_{\infty}>2L_b^{-1}(lpha)
ight)\leq lpha+O\left(\sqrt{1/n}
ight)$$

Subsampling

A Variant of the Bootstrap

Let S_b^i be a subsample of size b < n, for i = 1, ..., B. L_b is the CDF of $H(S_b^i, S_n)$.

Confidence Sets from Subsampling [FLRWBS]

Assume that p(x) is bounded away from zero. Then, almost surely, for all large n,

$$\mathbb{P}\left(||\textit{d}_{\mathbb{M}}-\textit{d}_{\mathcal{S}_n}||_{\infty}>2L_b^{-1}(lpha)
ight)\leq lpha+O\left(\sqrt{1/n}
ight)$$

[RS-2002] On the uniform asymptotic validity of subsampling and the bootstrap. Annals of Statistics.

Putting It All Together

Subsampling Theorem

 $\mathbb{P}\left(||\boldsymbol{d}_{\mathbb{M}}-\boldsymbol{d}_{\mathcal{S}_n}||_{\infty}>2L_b^{-1}(\alpha)\right)\leq \alpha+O\left(\sqrt{1/n}\right)$

Putting It All Together

Subsampling Theorem

$$\mathbb{P}\left(||\textit{d}_{\mathbb{M}}-\textit{d}_{\mathcal{S}_n}||_{\infty}>2L_b^{-1}(lpha)
ight)\leq lpha+O\left(\sqrt{1/n}
ight)$$

Bottleneck Stability Theorem

$$||d_{\mathbb{M}} - d_{\mathcal{S}_n}||_{\infty} \geq W_{\infty}(D,\widehat{D}_n)$$

Distance Function

Putting It All Together

Subsampling Theorem

$$\mathbb{P}\left(||\textit{d}_{\mathbb{M}}-\textit{d}_{\mathcal{S}_n}||_{\infty}>2L_b^{-1}(lpha)
ight)\leq lpha+O\left(\sqrt{1/n}
ight)$$

Bottleneck Stability Theorem

$$||d_{\mathbb{M}} - d_{\mathcal{S}_n}||_{\infty} \geq W_{\infty}(D, \widehat{D}_n)$$

Confidence Sets for Persistence Diagrams

$$\mathbb{P}\left(W_{\infty}(D,\widehat{D}_{n})>2L_{b}^{-1}(\alpha)\right)\leq\alpha+O\left(\sqrt{1/n}\right)$$

Distance Function

Putting It All Together

Subsampling Theorem

$$\mathbb{P}\left(||\textit{d}_{\mathbb{M}}-\textit{d}_{\mathcal{S}_n}||_{\infty}>2L_b^{-1}(lpha)
ight)\leq lpha+O\left(\sqrt{1/n}
ight)$$

Bottleneck Stability Theorem

$$||d_{\mathbb{M}} - d_{\mathcal{S}_n}||_{\infty} \geq W_{\infty}(D, \widehat{D}_n)$$

Confidence Sets for Persistence Diagrams

$$\mathbb{P}\left(W_{\infty}(D,\widehat{D}_{n})>2L_{b}^{-1}(\alpha)
ight)\leq lpha+O\left(\sqrt{1/n}
ight)$$

Asymptotic Confidence Sets for Persistence Diagrams

$$\lim_{n\to\infty} \mathbb{P}\left(W_{\infty}(D,\widehat{D}_n) \leq 2L_b^{-1}(\alpha)\right) \geq 1-\alpha$$

B. Fasy (Tulane)

Statistics and Topology

Varying α



Varying α



Two More Methods

 $\mathcal{S}_n = \mathcal{S}_{1,n} \bigsqcup \mathcal{S}_{2,n}.$

Theorem (Concentration of Measure)

There exists $\hat{t}_{\textit{cm}} = \hat{t}_{\textit{cm}}(\alpha, d, n, \mathcal{S}_{1,n})$ such that

$$\mathbb{P}\left(W_{\infty}(D,\widehat{D}_n) > \widehat{t}_{cm}\right) \leq \alpha + O\left(\left(\frac{\log n}{n}\right)^{1/d+2}\right)$$

Theorem (Method of Shells)

There exists $\hat{t}_s = \hat{t}_s(\alpha, d, n, K, \mathcal{S}_{1,n})$ such that

$$\mathbb{P}\left(W_{\infty}(D,\widehat{D}_n) > \widehat{t}_s\right) \leq \alpha + O\left(\left(\frac{\log n}{n}\right)^{1/d+2}\right)$$

These Methods are Different



Distance Function Examples

Uniform Distribution on Unit Circle



Distance Function Examples

Uniform Distribution on Cassini Curve



Distance Function Examples

Cassini Curve with Outliers


Distance Function Examples

Normal Distribution on Unit Circle



The Intersection of Statistics and Topology: Stochastic Convergence of Persistence Landscapes and Silhouettes

Brittany Terese Fasy

joint work with F. Chazal, F. Lecci, B. Michel A. Rinaldo, L. Wasserman

3 June 2014

[CFLRW] Stochastic Convergence of Persistence Landscapes and Silhouettes. ArXiv. SoCG 2014 Proceedings, Kyoto. Journal version pending.









[B-2012] Statistical Topology Using Persistent Homology. ArXiv 1207.6437

B. Fasy (Tulane)

Statistics and Topology

3 June 2014 30 / 48

Weak Convergence of Landscapes



Weak Convergence of Landscapes



Weak Convergence of Landscapes



Let $\lambda_1, \ldots, \lambda_n \sim \mathcal{L}_T$. λ : true (unknown) landscape $\mu = \mathbb{E}(\lambda_i)$ $\bar{\lambda}_n$: average landscape

Weak Convergence of Landscapes



Pointwise Convergence [B-2012].

 $\bar{\lambda}_n$ converges pointwise to μ .

Weak Convergence of Landscapes



Pointwise Convergence [B-2012].

 $\bar{\lambda}_n$ converges pointwise to μ .

Gaussian Process on [0, T]

For $t \in [0, T]$, we define $\mathbb{G}_n(f_t) = \mathbb{G}_n(t) := \frac{1}{\sqrt{n}} (\overline{\lambda}_n(t) - \mu(t)).$

Uniform Convergence

Weak Convergence

 \mathbb{G}_n converges weakly to the Brownian bridge \mathbb{G} with covariance function

$$\kappa(f,g) = \int f(u)g(y)dP(u) - (\int f(u)dP(u))(\int g(u)dP(u)).$$

Uniform Convergence

Weak Convergence

 \mathbb{G}_n converges weakly to the Brownian bridge \mathbb{G} with covariance function

$$\kappa(f,g) = \int f(u)g(y)dP(u) - (\int f(u)dP(u))(\int g(u)dP(u)).$$

Uniform CLT

There exists a random variable $W \stackrel{d}{=} \sup_{t \in [t_*, t^*]} |\mathbb{G}(f_t)|$ such that

$$\sup_{z\in\mathbb{R}}\left|\mathbb{P}\Big(\sup_{t\in[t_*,t^*]}|\mathbb{G}_n(t)|\leq z\Big)-\mathbb{P}(W\leq z)\right|=O\Big(\frac{(\log n)^{\frac{7}{8}}}{n^{\frac{1}{8}}}\Big).$$

Uniform Convergence

Weak Convergence

 \mathbb{G}_n converges weakly to the Brownian bridge \mathbb{G} with covariance function

$$\kappa(f,g) = \int f(u)g(y)dP(u) - (\int f(u)dP(u))(\int g(u)dP(u)).$$

Uniform CLT

There exists a random variable $W \stackrel{d}{=} \sup_{t \in [t_*, t^*]} |\mathbb{G}(f_t)|$ such that

$$\sup_{z\in\mathbb{R}}\left|\mathbb{P}\Big(\sup_{t\in[t_*,t^*]}|\mathbb{G}_n(t)|\leq z\Big)-\mathbb{P}(W\leq z)\right|=O\Big(\frac{(\log n)^{\frac{7}{8}}}{n^{\frac{1}{8}}}\Big).$$

Proofs rely on [CCK-2013]: Anti-concentration and honest adaptive confidence bands. ArXiv 1303.7152.

B. Fasy (Tulane)

Confidence Bands



Confidence Bands



Asymptotic Confidence Bands



$$\begin{split} \xi_{1}, \dots, \xi_{n} \sim \mathcal{N}(0, 1) \\ & \widetilde{\mathbb{G}}_{n}(f_{t}) \coloneqq \frac{1}{\sqrt{n}} \sum_{i=1}^{n} \xi_{i}(\lambda_{i}(t) - \overline{\lambda}_{n}(t)) \\ \\ & \textbf{\alpha-Quantile} \\ & \widetilde{Z}_{\alpha} \text{ is the unique value such that} \\ & \mathbb{P}\left(\sup_{t} |\widetilde{\mathbb{G}}_{n}(f_{t})| > \widetilde{Z}_{\alpha} \ \Big| \ \{\lambda_{i}\}\right) = \alpha \end{split}$$



The Multiplier Bootstrap



Uniform Band

$$\mathbb{P}\left(\ell_n(t) \leq \mu(t) \leq u_n(t) \text{ for all } t
ight) \geq 1 - lpha - \Big(rac{\left(\log n
ight)^{rac{l}{8}}}{n^{rac{1}{8}}}\Big).$$

Variable Width Confidence Bands

$$\mathbb{H}_n(f_t) := \mathbb{G}_n(t) / \sigma(t) = \frac{1}{\sqrt{n}} \sum_{i=1}^n \frac{\lambda_i(t) - \mu(t)}{\sigma(t)}$$

$$\widehat{\mathbb{H}}_n(f_t) := \frac{1}{\sqrt{n}} \sum_{i=1}^n \xi_i \frac{\lambda_i(t) - \overline{\lambda}_n(t)}{\widehat{\sigma}_n(t)}$$

 \widehat{Q}_{lpha} is the unique value such that

$$\mathbb{P}\left(\sup_{t} |\widetilde{\mathbb{H}}_{n}(f_{t})| > \widetilde{Q}_{\alpha} \mid \{\lambda_{i}\}\right) = \alpha$$

The Multiplier Bootstrap



The Intersection of Statistics and Topology: ... And More!

Brittany Terese Fasy

joint work with S. Balakrishnan, F. Lecci, A. Rinaldo, A. Singh, L. Wasserman

3 June 2014

Definitions



Definitions



Definitions



Definitions



Weighted Silhouette

$$\phi(t) = \frac{\sum_{i=1}^{n} w_i \Lambda_j(t)}{\sum_{j=1}^{m} \sum w_j}$$

Power-Weighted Silhouette

$$w_i = |d_i - b_i|^p$$

Power-Weighted Silhouettes

Two Examples



Results

Since ϕ is one-Lipschitz for non-negative weights $w_j \dots$

Convergence of Empirical Process

$$\frac{1}{\sqrt{n}}\left(\sum_{i=1}^n \phi_i(t) - \mathbb{E}[\phi(t)]\right)$$

converges weakly to a Brownian bridge, with known rate of convergence.

Confidence Bands

We can use the multiplier bootstrap to create a uniform (or a variable width) confidence band defined by ℓ_n^{sil} and u_n^{sil} such that

$$\lim_{n\to\infty}\mathbb{P}\left(\ell_n^{\textit{sil}}(t)\leq \mu(t)\leq u_n^{\textit{sil}}(t) \text{ for all } t\right)=1-\alpha.$$

Example I A Toy Example



Example II Earthquake Epicenters



B. Fasy (Tulane)

Statistics and Topology

Summary



Collaborator Collage



Thank you!



References

[CDGO] The Structure and Stability of Persistence Modules. ArXiv 1207.3674.

[CFLRSW] On the Bootstrap for Persistence Diagrams and Landscapes. Modeling and Analysis of Information Systems, **20**:6 (Dec. 2013), 96–105.



