



UCL



GATSBY

Inria

KSD Aggregated Goodness-of-fit Tests KSDAgg & KSDAggInc

Antonin Schrab

University College London
Centre for Artificial Intelligence
Gatsby Computational Neuroscience Unit
Inria London

a.schrab@ucl.ac.uk

antoninschrab.github.io

KSD Aggregated Goodness-of-fit Test



Antonin
Schrab

†‡§



Benjamin
Guedj

†§



Arthur
Gretton

‡

† Centre for Artificial Intelligence, UCL

‡ Gatsby Computational Neuroscience Unit, UCL

§ Inria London

Efficient Aggregated Kernel Tests using Incomplete U -statistics



Antonin
Schrab

†‡§



Ilmun
Kim

*



Benjamin
Guedj

†§



Arthur
Gretton

‡

† Centre for Artificial Intelligence, UCL

‡ Gatsby Computational Neuroscience Unit, UCL

§ Inria London

* Department of Statistics & Data Science, Yonsei University

Kernel Stein Discrepancy

Goodness-of-fit problem:

- model with probability density p or score function $\nabla \log p(\cdot)$ on \mathbb{R}^d
- samples $Z_N := (Z_1, \dots, Z_N)$, $Z_i \stackrel{\text{iid}}{\sim} q$ in \mathbb{R}^d

$$\mathcal{H}_0: p = q \quad \text{against} \quad \mathcal{H}_a: p \neq q$$

Kernel Stein Discrepancy

Goodness-of-fit problem:

- model with probability density p or score function $\nabla \log p(\cdot)$ on \mathbb{R}^d
- samples $Z_N := (Z_1, \dots, Z_N)$, $Z_i \stackrel{\text{iid}}{\sim} q$ in \mathbb{R}^d

$$\mathcal{H}_0: p = q \quad \text{against} \quad \mathcal{H}_a: p \neq q$$

Stein kernel: $h_{p,k}(x, y)$ defined in terms of $\nabla \log p(\cdot)$ and kernel k

Kernel Stein Discrepancy

Goodness-of-fit problem:

- model with probability density p or score function $\nabla \log p(\cdot)$ on \mathbb{R}^d
- samples $Z_N := (Z_1, \dots, Z_N)$, $Z_i \stackrel{\text{iid}}{\sim} q$ in \mathbb{R}^d

$$\mathcal{H}_0: p = q \quad \text{against} \quad \mathcal{H}_a: p \neq q$$

Stein kernel: $h_{p,k}(x, y)$ defined in terms of $\nabla \log p(\cdot)$ and kernel k

Stein identity: $\mathbb{E}_p[h_{p,k}(Z, \cdot)] = 0$

Kernel Stein Discrepancy

Goodness-of-fit problem:

- model with probability density p or score function $\nabla \log p(\cdot)$ on \mathbb{R}^d
- samples $Z_N := (Z_1, \dots, Z_N)$, $Z_i \stackrel{\text{iid}}{\sim} q$ in \mathbb{R}^d

$$\mathcal{H}_0: p = q \quad \text{against} \quad \mathcal{H}_a: p \neq q$$

Stein kernel: $h_{p,k}(x, y)$ defined in terms of $\nabla \log p(\cdot)$ and kernel k

Stein identity: $\mathbb{E}_{p,k}[h_{p,k}(Z, \cdot)] = 0$

KSD: $\text{KSD}_{p,k}^2(q) := \mathbb{E}_{q,q}[h_{p,k}(Z, Z')]$

Kernel Stein Discrepancy

Goodness-of-fit problem:

- model with probability density p or score function $\nabla \log p(\cdot)$ on \mathbb{R}^d
- samples $\mathbb{Z}_N := (Z_1, \dots, Z_N)$, $Z_i \stackrel{\text{iid}}{\sim} q$ in \mathbb{R}^d

$$\mathcal{H}_0: p = q \quad \text{against} \quad \mathcal{H}_a: p \neq q$$

Stein kernel: $h_{p,k}(x, y)$ defined in terms of $\nabla \log p(\cdot)$ and kernel k

Stein identity: $\mathbb{E}_{p,k}[h_{p,k}(Z, \cdot)] = 0$

KSD: $\text{KSD}_{p,k}^2(q) := \mathbb{E}_{q,q}[h_{p,k}(Z, Z')]$

U-statistic: $\widehat{\text{KSD}}_{p,k}^2(\mathbb{Z}_N) := \frac{1}{N(N-1)} \sum_{1 \leq i \neq j \leq N} h_{p,k}(Z_i, Z_j)$

Kernel Stein Discrepancy

Goodness-of-fit problem:

- model with probability density p or score function $\nabla \log p(\cdot)$ on \mathbb{R}^d
- samples $\mathbb{Z}_N := (\mathcal{Z}_1, \dots, \mathcal{Z}_N)$, $\mathcal{Z}_i \stackrel{\text{iid}}{\sim} q$ in \mathbb{R}^d

$$\mathcal{H}_0: p = q \quad \text{against} \quad \mathcal{H}_a: p \neq q$$

Stein kernel: $h_{p,k}(x, y)$ defined in terms of $\nabla \log p(\cdot)$ and kernel k

Stein identity: $\mathbb{E}_{p,k}[h_{p,k}(\mathcal{Z}, \cdot)] = 0$

KSD: $\text{KSD}_{p,k}^2(q) := \mathbb{E}_{q,q}[h_{p,k}(\mathcal{Z}, \mathcal{Z}')]$

U-statistic: $\widehat{\text{KSD}}_{p,k}^2(\mathbb{Z}_N) := \frac{1}{N(N-1)} \sum_{1 \leq i \neq j \leq N} h_{p,k}(\mathcal{Z}_i, \mathcal{Z}_j)$

Incomplete U-statistic: $\overline{\text{KSD}}_{p,k}^2(\mathbb{Z}_N) := \frac{1}{|\mathcal{D}|} \sum_{(i,j) \in \mathcal{D}} h_{p,k}(\mathcal{Z}_i, \mathcal{Z}_j)$

KSD Goodness-of-fit Tests

KSD test: reject $\mathcal{H}_0: p = q$ if $\widehat{\text{KSD}}_{p,k}^2(\mathbb{Z}_N) > \widehat{q}_{1-\alpha}^k$

- quantile $\widehat{q}_{1-\alpha}^k$ estimated using wild/parametric bootstrap
- complexity $\mathcal{O}(B N^2)$

KSD Goodness-of-fit Tests

KSD test: reject $\mathcal{H}_0: p = q$ if $\widehat{\text{KSD}}_{p,k}^2(\mathbb{Z}_N) > \widehat{q}_{1-\alpha}^k$

- quantile $\widehat{q}_{1-\alpha}^k$ estimated using wild/parametric bootstrap
- complexity $\mathcal{O}(B N^2)$

KSDAgg: reject \mathcal{H}_0 if $\widehat{\text{KSD}}_{p,k}^2(\mathbb{Z}_N) > \widehat{q}_{1-u_\alpha w_k}^k$ for some $k \in \mathcal{K}$

- positive weights $(w_k)_{k \in \mathcal{K}}$ summing to 1
- correction u_α maximizes power while retaining well-calibrate level α
- complexity $\mathcal{O}(|\mathcal{K}| B N^2)$

KSD Goodness-of-fit Tests

KSD test: reject $\mathcal{H}_0: p = q$ if $\widehat{\text{KSD}}_{p,k}^2(\mathbb{Z}_N) > \widehat{q}_{1-\alpha}^k$

- quantile $\widehat{q}_{1-\alpha}^k$ estimated using wild/parametric bootstrap
- complexity $\mathcal{O}(B N^2)$

KSDAgg: reject \mathcal{H}_0 if $\widehat{\text{KSD}}_{p,k}^2(\mathbb{Z}_N) > \widehat{q}_{1-u_\alpha w_k}^k$ for some $k \in \mathcal{K}$

- positive weights $(w_k)_{k \in \mathcal{K}}$ summing to 1
- correction u_α maximizes power while retaining well-calibrate level α
- complexity $\mathcal{O}(|\mathcal{K}| B N^2)$

KSDAggInc: reject \mathcal{H}_0 if $\overline{\text{KSD}}_{p,k}^2(\mathbb{Z}_N) > \widehat{q}_{1-u_\alpha w_k}^k$ for some $k \in \mathcal{K}$

- complexity $\mathcal{O}(|\mathcal{K}| B |\mathcal{D}|)$
- linear-time test for the choice $|\mathcal{D}| = c N$ for some small $c \in \mathbb{N}$

- **Power guarantees:** upper bound on uniform separation rates
- **Trade-off:** computational efficiency versus rate of convergence

Come check out my poster for details and discussions!

- **Power guarantees:** upper bound on uniform separation rates
- **Trade-off:** computational efficiency versus rate of convergence

KSD Aggregated Goodness-of-fit Tests

KSDAgg: KSD Aggregated Goodness-of-fit Test
KSDAggInc: Efficient Aggregated Kernel Tests using Incomplete U-statistics

Contributions

- Aggregate KSD tests with different kernels or bandwidths
- Quantiles estimated via wild or parametric bootstraps
- No data splitting (known to result in a loss in power)
- Uniform separation rate upper bound for general kernels
- Propose efficient tests based on incomplete U-statistics
- Quantify trade-off efficiency versus rate of convergence

Goodness-of-fit problem

Are samples drawn from the model?

- model density p (or score function $\nabla \log p(x)$)
- samples $\mathbf{Z}_N := (\mathbf{z}_1, \dots, \mathbf{z}_N)$ drawn $\mathbf{z}_i \stackrel{\text{iid}}{\sim} p$

Hypothesis testing:

$$H_0: p = q \quad \text{against} \quad H_a: p \neq q$$

Kernel Stein Discrepancy

Stein kernel: $h_{p,k}(x, y)$ in terms of $\nabla \log p(z)$ with kernel k

Stein Identity: $\mathbb{E}_p[h_p(\cdot, \mathbf{Z}_N)] = 0$

Kernel Stein Discrepancy: $\text{KSD}_{p,k}^2(q) := \mathbb{E}_q[h_q(\mathbf{Z}_N, \mathbf{Z}')]^2$

Estimator: $\widehat{\text{KSD}}_{p,k}^2(\mathbf{Z}_N) := \frac{1}{N(N-1)} \sum_{1 \leq i < j \leq N} h_{p,k}(\mathbf{z}_i, \mathbf{z}_j)$

KSD test for fixed kernel k

Test: reject H_0 if $\widehat{\text{KSD}}_{p,k}^2(\mathbf{Z}_N) > \widehat{q}_{1-\alpha}^k$

Quantile: $\widehat{q}_{1-\alpha}^k = (1 - \alpha)^{-1}\text{-th}$ largest bootstrapped value

Wild bootstrap: $\frac{1}{N(N-1)} \sum_{1 \leq i < j \leq N} h_{p,k}(\mathbf{z}_i, \mathbf{z}_j), \mathbf{z}_i \stackrel{\text{iid}}{\sim} p$

Parametric bootstrap: $\frac{1}{N(N-1)} \sum_{1 \leq i < j \leq N} h_{p,k}(\tilde{\mathbf{z}}_i, \tilde{\mathbf{z}}_j), \tilde{\mathbf{z}}_i \stackrel{\text{iid}}{\sim} p$

KSDAgg for collection of kernels \mathcal{K}

Test: reject H_0 if $\widehat{\text{KSD}}_{p,\mathcal{K}}^2(\mathbf{Z}_N) > \widehat{q}_{1-\alpha}^{\mathcal{K}}$ for some $k \in \mathcal{K}$

Weights (prior): $(w_k)_{k \in \mathcal{K}}$ satisfying $\sum_{k \in \mathcal{K}} w_k \leq 1$

Correction: u_0 maximum value such that the level estimated via Monte-Carlo is well-calibrated at α

More powerful than conservative Bonferroni correction

KSDAgg Uniform separation rate

Integral transform: $(T_\alpha f)(y) := \int_{\mathbb{R}^d} \kappa(x, y) f(x) dx$

Kernel assumption: $A_k := \mathbb{E}_q[h_k(\mathbf{Z}_N, \mathbf{Z}')^2] < \infty$

If $\|p - q\|_2^2$ is greater than

$$\min_{1 \leq k \leq K} \left(\|p - q\| - T_{h_k}(p - q) \right)^2 + C N^{-1} \ln \left(\frac{1}{\alpha u_0} \right) \sqrt{\frac{A_k}{\beta}}$$

then KSDAgg has power at least $1 - \beta$.

Incomplete U-statistic

Estimator: $\widehat{\text{KSD}}_{p,\mathcal{K}}^2(\mathbf{Z}_N) := \frac{1}{N(N-1)} \sum_{(i,j) \in \mathcal{D}_N} h_{p,\mathcal{K}}(\mathbf{z}_i, \mathbf{z}_j)$

Design: \mathcal{D}_N random / deterministic subset of $\{(i, j)\}_{1 \leq i < j \leq N}$

Linear time: $|\mathcal{D}_N| = cN$ for some fixed integer $c \in \mathbb{N}$

KSDAggInc Uniform separation rate

KSDAggInc: use $\widehat{\text{KSD}}_{p,\mathcal{K}}^2(\mathbf{Z}_N)$ instead of $\widehat{\text{KSD}}_{p,\mathcal{K}}^2(\mathbf{Z}_N)$

Uniform separation rate: same condition as for KSDAgg with N multiplied by an extra cost factor $|\mathcal{D}_N|/N^2$

- $|\mathcal{D}_N| \propto N^2$: recover KSDAgg rate
- $N \lesssim |\mathcal{D}_N| \lesssim N^2$: cost $|\mathcal{D}_N|/N^2$ incurred in KSDAgg rate
- Trade-off: computational efficiency / rate convergence
- $|\mathcal{D}_N| \lesssim N$: no guarantee that rate converges to 0

Experiments

Gaussian-Bernoulli Restricted Boltzmann Machine: graphical model with binary hidden variable $h \in \{\pm 1\}^{d_h}$ & continuous observable variable $x \in \mathbb{R}^{d_x}$ with joint density

$$p(x, h) = \frac{1}{Z} \exp \left(\frac{1}{2} x^\top B h + B^\top x + c^\top h - \frac{1}{2} \|x\|_2^2 \right)$$

- model: GBRBM with $B \in \{\pm 1\}^{d_x \times d_h}$, $b \in \mathbb{R}^{d_h}$, $c \in \mathbb{R}^{d_x}$
- samples: GBRBM with noise $\mathcal{N}(0, \sigma)$ injected into B

Collection: Gaussian kernels with scaled median bandwidth

Parameter R : number of subdiagonals of the kernel matrix

FSSD: Jitkrittum et al. 2017 LSD: Grathwohl et al. 2020
L1 IQM & Cauchy RFF: Huggins and Mackey 2018

Power

Sample size N

$d_h = 50, d_x = 40, \sigma = 0.02$

Time (seconds)

Sample size N

$d_h = 50, d_x = 40, \sigma = 0.02$

Power

Dimension d_h

$d_x = 100, N = 1000, \sigma = 0.01$

Power

Noise standard deviation σ

$d_h = 50, d_x = 40, \sigma = 0.02$

KSDAgg

paper

code

AgInnc

paper

code

Antonin Schrab

KSDAgg & KSDAggInc

4 / 4