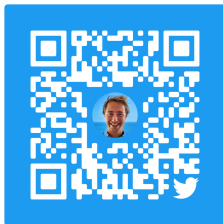


## Aggregated Kernel Tests



**Antonin Schrab**

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

[a.schrab@ucl.ac.uk](mailto:a.schrab@ucl.ac.uk)

[antoninschrab.github.io](https://antoninschrab.github.io)

- 1 MMDAgg: MMD Aggregated Two-Sample Test
- 2 KSDAgg: KSD Aggregated Goodness-of-fit Test
- 3 AggInc: Efficient Aggregated Kernel Tests using Incomplete  $U$ -statistics

# MMD Aggregated Two-Sample Test



Antonin  
Schrab

† ‡ §



Ilmun  
Kim

\*



Mélisande  
Albert

\*



Béatrice  
Laurent

\*



Benjamin  
Guedj

† §



Arthur  
Gretton

‡

† Centre for Artificial Intelligence, UCL

‡ Gatsby Computational Neuroscience Unit, UCL

§ Inria London Programme

\* Department of Statistics & Data Science, Yonsei University

\* Institut de Mathématiques, Université de Toulouse

# Two-sample problem

- samples  $\mathbb{X}_m := (X_1, \dots, X_m)$ ,  $X_i \stackrel{\text{iid}}{\sim} p$  in  $\mathbb{R}^d$
- samples  $\mathbb{Y}_n := (Y_1, \dots, Y_n)$ ,  $Y_i \stackrel{\text{iid}}{\sim} q$  in  $\mathbb{R}^d$

$$\begin{array}{lll} \mathcal{H}_0: p = q & \text{against} & \mathcal{H}_a: p \neq q \\ \Delta(\mathbb{X}_m, \mathbb{Y}_n) = 1 & \iff & \text{reject } \mathcal{H}_0 \\ \Delta(\mathbb{X}_m, \mathbb{Y}_n) = 0 & \iff & \text{fail to reject } \mathcal{H}_0 \end{array}$$

**Type I error:** controlled by  $\alpha$  by design

$$\mathbb{P}_{p \times p}(\Delta(\mathbb{X}_m, \mathbb{Y}_n) = 1) \leq \alpha$$

**Type II error:** find a condition on  $\|p - q\|_2$  to control by  $\beta$

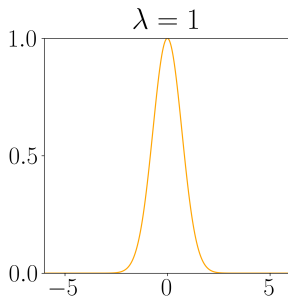
$$\mathbb{P}_{p \times q}(\Delta(\mathbb{X}_m, \mathbb{Y}_n) = 0) \leq \beta$$



**Kernel:**  $k_\lambda(\mathbf{x}, \mathbf{y}) := \prod_{i=1}^d \frac{1}{\lambda_i} K_i\left(\frac{x_i - y_i}{\lambda_i}\right)$     **Bandwidth:**  $\lambda \in (0, \infty)^d$

**Gaussian kernel:**  $K_i(u) = \frac{1}{\sqrt{\pi}} \exp(-u^2)$ ,  $u \in \mathbb{R}$ ,  $i = 1, \dots, d$

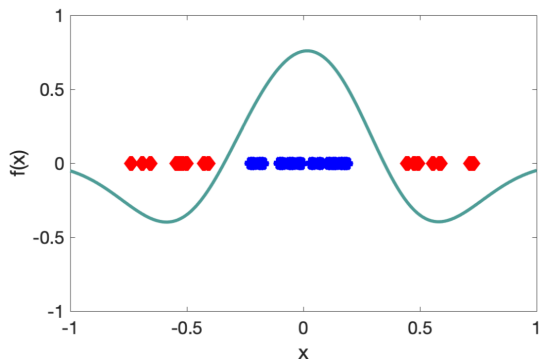
$$k_\lambda(\mathbf{x}, \mathbf{y}) := \frac{1}{\pi^{d/2} \lambda_1 \dots \lambda_d} \exp\left(-\sum_{i=1}^d \frac{(x_i - y_i)^2}{\lambda_i^2}\right)$$



# Two-sample test using the Maximum Mean Discrepancy

**Kernel:**  $k_\lambda(x, y) := \prod_{i=1}^d \frac{1}{\lambda_i} \mathcal{K}_i\left(\frac{x_i - y_i}{\lambda_i}\right)$       **Bandwidth:**  $\lambda \in (0, \infty)^d$

$$\text{MMD}_\lambda(p, q) := \sup_{f \in \mathcal{H}_\lambda: \|f\|_{\mathcal{H}_\lambda} \leq 1} \left| \mathbb{E}_{X \sim p}[f(X)] - \mathbb{E}_{Y \sim q}[f(Y)] \right|$$



$p \neq q$

# Bandwidth intuition

- **Small sample sizes:** only global differences are detectable
  - **Small bandwidth:** wrongly detects artificial local differences under  $\mathcal{H}_0$
  - **Large bandwidth:** well-suited to detect global differences under  $\mathcal{H}_a$
- **Large sample sizes:** local differences are detectable
  - **Small bandwidth:** well-suited to detect local differences under  $\mathcal{H}_a$
  - **Large bandwidth:** fails to detect local differences under  $\mathcal{H}_a$

⇒ **Bandwidths** should decrease as the **sample sizes** increase

- Choice of **bandwidth** is **crucial** for test power!
- **Bandwidth** selection methods: **median heuristic** & **data splitting**
- **Our method:** aggregate multiple tests with different **bandwidths**

# Maximum Mean Discrepancy estimator

$$\begin{aligned} \text{MMD}_\lambda^2(p, q) &:= \mathbb{E}_{p,p}[k_\lambda(X, X')] \\ &\quad - 2\mathbb{E}_{p,q}[k_\lambda(X, Y)] \\ &\quad + \mathbb{E}_{q,q}[k_\lambda(Y, Y')] \end{aligned}$$

$$\begin{aligned} \widehat{\text{MMD}}_\lambda^2(\mathbb{X}_m, \mathbb{Y}_n) &:= \frac{1}{m(m-1)} \sum_{1 \leq i \neq i' \leq m} k_\lambda(X_i, X_{i'}) \\ &\quad - \frac{2}{mn} \sum_{i=1}^m \sum_{j=1}^n k_\lambda(X_i, Y_j) \\ &\quad + \frac{1}{n(n-1)} \sum_{1 \leq j \neq j' \leq n} k_\lambda(Y_j, Y_{j'}) \end{aligned}$$

## MMD test for a fixed bandwidth $\lambda$

$$\Delta_\alpha^\lambda(\mathbb{X}_m, \mathbb{Y}_n) := \mathbb{1}\left(\widehat{\text{MMD}}_\lambda^2(\mathbb{X}_m, \mathbb{Y}_n) > \widehat{q}_{1-\alpha}^\lambda\right)$$

**Quantile:**  $\widehat{q}_{1-\alpha}^\lambda$  is the  $[(B+1)(1-\alpha)]$ -th largest value of  $\widehat{\text{MMD}}_\lambda^2(\mathbb{X}_m, \mathbb{Y}_n)$  and  $B$   $\mathcal{H}_0$ -simulated test statistics

**Permutations:**  $\widehat{\text{MMD}}_\lambda^2(\mathbb{X}_m^\sigma, \mathbb{Y}_n^\sigma)$  where  $(\mathbb{X}_m^\sigma, \mathbb{Y}_n^\sigma) = \sigma(\mathbb{X}_m \cup \mathbb{Y}_n)$

**Wild bootstrap:** case  $m = n$ ,  $\epsilon_1, \dots, \epsilon_n \stackrel{\text{iid}}{\sim} \text{Unif}\{-1, 1\}$  (Rademacher)

$$\frac{1}{n(n-1)} \sum_{1 \leq i \neq j \leq n} \epsilon_i \epsilon_j \left( k_\lambda(X_i, X_j) - k_\lambda(X_i, Y_j) - k_\lambda(Y_i, X_j) + k_\lambda(Y_i, Y_j) \right)$$

**Non-asymptotic level:**  $\mathbb{P}_{p \times p}(\Delta(\mathbb{X}_m, \mathbb{Y}_n) = 1) \leq \alpha$

**Time complexity:**  $\mathcal{O}(B(m+n)^2)$

# MMDAgg for a collection of bandwidths $\Lambda$

**Bonferroni multiple testing:** non-asymptotic level  $\alpha$

$$\Delta_{\alpha}^{\Lambda}(\mathbb{X}_m, \mathbb{Y}_n) := \mathbb{1} \left( \widehat{\text{MMD}}_{\lambda}^2(\mathbb{X}_m, \mathbb{Y}_n) > \hat{q}_{1-\alpha/|\Lambda|}^{\lambda} \text{ for some } \lambda \in \Lambda \right)$$

- time complexity  $\mathcal{O}(|\Lambda| B_1 (m+n)^2)$

**MMDAgg:** non-asymptotic level  $\alpha$

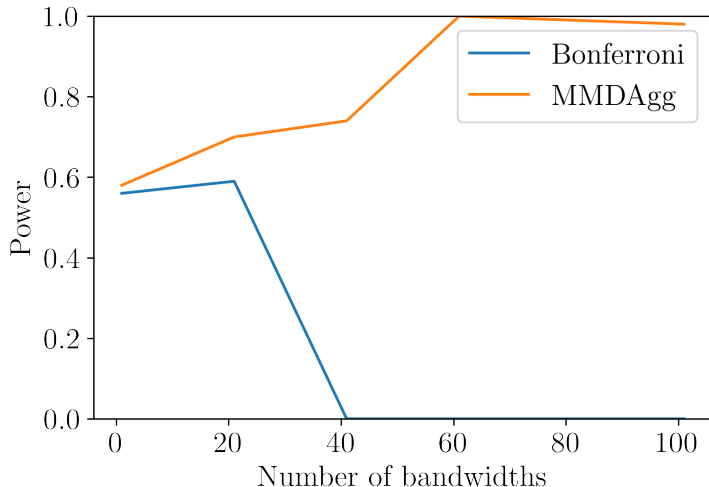
$$\Delta_{\alpha}^{\Lambda}(\mathbb{X}_m, \mathbb{Y}_n) := \mathbb{1} \left( \widehat{\text{MMD}}_{\lambda}^2(\mathbb{X}_m, \mathbb{Y}_n) > \hat{q}_{1-u_{\alpha} w_{\lambda}}^{\lambda} \text{ for some } \lambda \in \Lambda \right)$$

- positive weights  $(w_{\lambda})_{\lambda \in \Lambda}$  satisfying  $\sum_{\lambda \in \Lambda} w_{\lambda} \leq 1$
- correction  $u_{\alpha}$  defined as

$$\sup \left\{ u > 0 : \mathbb{P}_{p \times p} \left( \max_{\lambda \in \Lambda} \left( \widehat{\text{MMD}}_{\lambda}^2(\mathbb{X}_m, \mathbb{Y}_n) - \hat{q}_{1-uw_{\lambda}}^{\lambda} \right) > 0 \right) \leq \alpha \right\}$$

- more powerful than Bonferroni correction as  $u_{\alpha} \geq \alpha$
- time complexity  $\mathcal{O}(|\Lambda| (B_1 + B_2) (m+n)^2)$

# Multiple testing correction comparison



$$\Lambda(\ell_-, \ell_+) := \{2^\ell \lambda_{med} : \ell \in \{\ell_-, \dots, \ell_+\}\} \quad w_\lambda := 1 / |\Lambda|$$
$$\Lambda(-i, i), \quad i \in \{0, 10, 20, 30, 40, 50\}$$

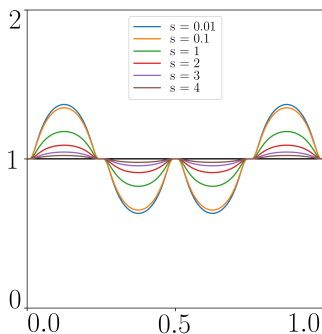
# Sobolev balls

Regularity/smoothness assumption:  $p - q \in \mathcal{S}_d^s(R)$

Sobolev balls:

$$\mathcal{S}_d^s(R) := \left\{ f \in L^1(\mathbb{R}^d) \cap L^2(\mathbb{R}^d) : \int_{\mathbb{R}^d} \|\xi\|_2^{2s} |\widehat{f}(\xi)|^2 d\xi \leq (2\pi)^d R^2 \right\}$$

- radius  $R > 0$
- dimension  $d$
- smoothness parameter  $s > 0$  (unknown)
- Fourier transform  $\widehat{f}(\xi) := \int_{\mathbb{R}^d} f(x) e^{-ix^\top \xi} dx$





## Theorem

$$\Lambda^* := \left\{ 2^{-\ell} \mathbb{1}_d : \ell \in \left\{ 1, \dots, \left\lceil \frac{2}{d} \log_2 \left( \frac{m+n}{\ln(\ln(m+n))} \right) \right\rceil \right\} \right\}, \quad w_\lambda := \frac{6}{\pi^2 \ell^2}$$

Assuming  $p - q \in \mathcal{S}_d^s(R)$ , the condition

$$\|p - q\|_2 \geq C \left( \frac{m+n}{\ln(\ln(m+n))} \right)^{-2s/(4s+d)}$$

guarantees control over the probability of type II error of MMDAgg

$$\mathbb{P}_{p \times q}(\Delta_\alpha^{\Lambda^*}(\mathbb{X}_m, \mathbb{Y}_n) = 0) \leq \beta.$$

**Minimax rate over Sobolev balls:**  $(m+n)^{-2s/(4s+d)}$

**Adaptive over**  $\{\mathcal{S}_d^s(R) : s > 0, R > 0\}$

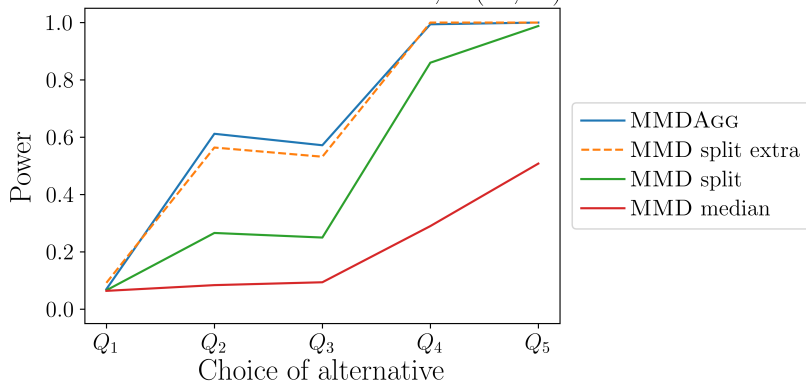
# MMDAgg Experiment

$$\Lambda(l_-, l_+) := \{2^l \lambda_{med} : l \in \{l_-, \dots, l_+\}\} \quad w_\lambda := 1 / |\Lambda|$$

$$\begin{aligned} \mathcal{P} &:= \{0, \dots, 9\} & \mathcal{Q}_2 &:= \mathcal{P} \setminus \{8, 6\} & \mathcal{Q}_4 &:= \mathcal{P} \setminus \{8, 6, 4, 2\} \\ \mathcal{Q}_1 &:= \mathcal{P} \setminus \{8\} & \mathcal{Q}_3 &:= \mathcal{P} \setminus \{8, 6, 4\} & \mathcal{Q}_5 &:= \mathcal{P} \setminus \{8, 6, 4, 2, 0\} \end{aligned}$$

Two-sample experiment

MNIST dataset  $m = n = 500$ ,  $\Lambda(12, 16)$



# KSD Aggregated Goodness-of-fit Test



Antonin  
Schrab

† ‡ §



Benjamin  
Guedj

† §



Arthur  
Gretton

‡

† Centre for Artificial Intelligence, UCL

‡ Gatsby Computational Neuroscience Unit, UCL

§ Inria London Programme

# Goodness-of-fit problem & Kernel Stein Discrepancy

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

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

**Stein kernel:**  $h_{p,\lambda}(x, y)$  defined as

$$\begin{aligned} & (\nabla \log p(x)^\top \nabla \log p(y)) k_\lambda(x, y) + \nabla \log p(y)^\top \nabla_1 k_\lambda(x, y) \\ & + \nabla \log p(x)^\top \nabla_2 k_\lambda(x, y) + \sum_{1 \leq i \leq d} \frac{\partial}{\partial x_i \partial y_i} k_\lambda(x, y) \end{aligned}$$

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

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

**Estimator:**  $\widehat{\text{KSD}}_{p,\lambda}^2(\mathbf{Z}_n) := \frac{1}{n(n-1)} \sum_{1 \leq i \neq j \leq n} h_{p,\lambda}(\mathbf{Z}_i, \mathbf{Z}_j)$

# KSDAgg: KSD Aggregated test

**Wild bootstrap:** asymptotic level  $\alpha$

$$\frac{1}{n(n-1)} \sum_{1 \leq i \neq j \leq n} \epsilon_i \epsilon_j h_{p,\lambda}(Z_i, Z_j) \quad \text{where} \quad \epsilon_i \sim \text{Unif}\{-1, 1\}$$

**Parametric bootstrap:** non-asymptotic level  $\alpha$

$$\frac{1}{n(n-1)} \sum_{1 \leq i \neq j \leq n} h_{p,\lambda}(\tilde{Z}_i, \tilde{Z}_j) \quad \text{where} \quad \tilde{Z}_i \stackrel{\text{iid}}{\sim} p$$

**KSDAgg:**

$$\Delta_\alpha^\wedge(\mathbb{X}_m, \mathbb{Y}_n) := \mathbb{1} \left( \widehat{\text{KSD}}_{p,\lambda}^2(\mathbb{X}_m, \mathbb{Y}_n) > \hat{q}_{1-u_\alpha w_\lambda}^\lambda \text{ for some } \lambda \in \Lambda \right)$$

**Time complexity:**  $\mathcal{O}(|\Lambda| (B_1 + B_2) n^2)$

# KSDAgg Experiment

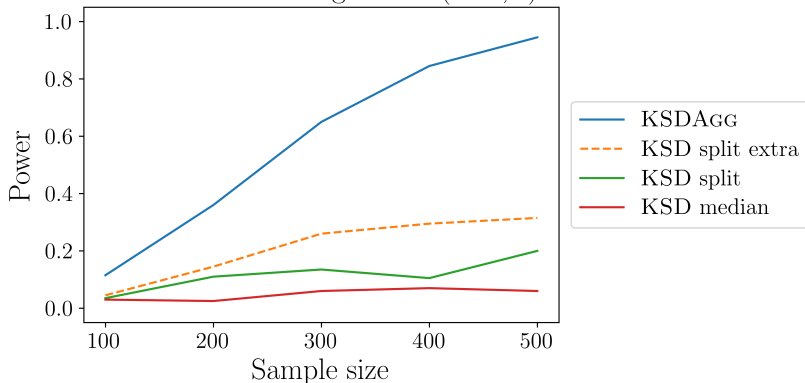
$$\Lambda(l_-, l_+) := \{2^l \lambda_{med} : l \in \{l_-, \dots, l_+\}\}$$

$$w_\lambda := 1 / |\Lambda|$$

**model:** Normalizing Flow density

**samples:** true MNIST digits

Goodness-of-fit experiment  
MNIST Normalizing Flow  $\Lambda(-20, 0)$



# What about HSICA<sub>agg</sub>?

## Independence problem:

Given paired samples  $((X_1, Y_1), \dots, (X_n, Y_n))$  in  $\mathbb{R}^{d_1} \times \mathbb{R}^{d_2}$  with

- joint probability density  $r$
- marginal probability densities  $p$  and  $q$

can we decide whether or not  $p \otimes q \neq r$  holds?

## Hilbert-Schmidt Independence Criterion:

$$\begin{aligned} \text{HSIC}_{k,\ell}(r) &:= \text{MMD}_{\kappa}(p \otimes q, r) \\ \kappa((X, Y), (X', Y')) &:= k(X, X')\ell(Y, Y') \end{aligned}$$

## ADAPTIVE TEST OF INDEPENDENCE BASED ON HSIC MEASURES.

Mélanie Albert<sup>\*1</sup>, Béatrice Laurent<sup>†1</sup>, Amandine Marrel<sup>‡2</sup>, and Anouar Meynaoui<sup>§1,2</sup>

<sup>1</sup>Institut de Mathématiques de Toulouse ; UMR5219, Université de Toulouse ; CNRS, INSA, F-31077 Toulouse, France.

<sup>2</sup>CEA, DEN, DER, F-13108 Saint-Paul-lez-Durance, France.

# 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 Programme

\* Department of Statistics & Data Science, Yonsei University



- **Complete  $U$ -statistic:**

$$\frac{1}{N(N-1)} \sum_{1 \leq i \neq j \leq N} h(Z_i, Z_j)$$

- **Quadratic-time MMDAgg, KSDAgg & HSICAgg:**

$$\mathcal{O}\left(|\Lambda| (B_1 + B_2) N^2\right)$$

- **Incomplete  $U$ -statistic:**

$$\frac{1}{|\mathcal{D}|} \sum_{(i,j) \in \mathcal{D}} h(Z_i, Z_j) \quad \mathcal{D} \subseteq \left\{ (i,j) : 1 \leq i \neq j \leq N \right\}$$

- **Efficient MMDAggInc, KSDAggInc & HSICAggInc:**

$$\mathcal{O}\left(|\Lambda| (B_1 + B_2) |\mathcal{D}|\right)$$

- **Linear-time** if  $|\mathcal{D}| = cN$

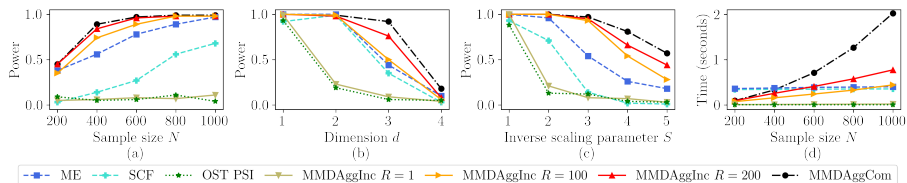
**Minimax rate over Sobolev balls:**  $N^{-2s/(4s+d)}$

**MMDAgg & HSICAgg rates:**  $\left( \frac{N}{\ln(\ln(N))} \right)^{-2s/(4s+d)}$

**MMDAggInc & HSICAggInc rates:**  $\left( \frac{|\mathcal{D}|/N}{\ln(\ln(|\mathcal{D}|/N))} \right)^{-2s/(4s+d)}$

- $|\mathcal{D}| \asymp N^2$ : recover the **Agg** rate
- $N \lesssim |\mathcal{D}| \lesssim N^2$ : cost  $|\mathcal{D}|/N^2$  incurred in the **Agg** rate  
**Trade-off:** computational efficiency / rate of convergence
- $|\mathcal{D}| \lesssim N$ : no guarantee that the **AggInc** rate converges to 0

## Perturbed uniform densities

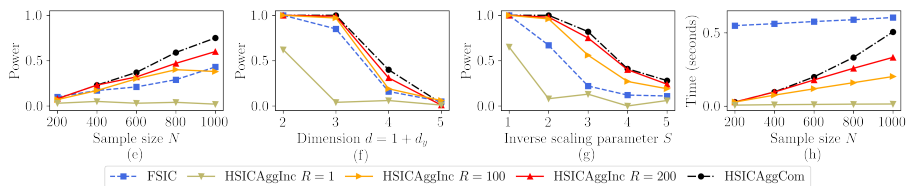


## ME (Mean Embeddings) & SCF (Smooth Characteristic Functions):

Jitkrittum, W., Szabó, Z., Chwialkowski, K. P., and Gretton, A. **Interpretable distribution features with maximum testing power**. In NeurIPS 2016.

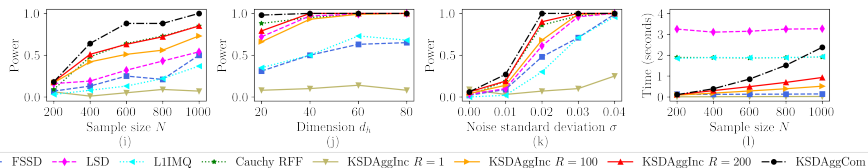
**OST PSI (One-sided Test Post Selection Inference):** Kübler, J. M., Jitkrittum, W., Schölkopf, B., and Muandet, K. **Learning kernel tests without data splitting**. NeurIPS 2020.

## Perturbed uniform densities



**FSIC (Finite Set Independence Criterion):** Jitkrittum, W., Szabó, Z., and Gretton, A. **An adaptive test of independence with analytic kernel embeddings.** In ICML 2017.

## Gaussian-Bernoulli Restricted Boltzmann Machine



**FSSD (Finite Set Stein Discrepancy):** Jitkrittum, W., Xu, W., Szabó, Z., Fukumizu, K., and Gretton, A. **A linear-time kernel goodness-of-fit test.** In NeurIPS 2017.

**LSD (Learned Stein Discrepancy):** Grathwohl, W., Wang, K.-C., Jacobsen, J.-H., Duvenaud, D., and Zemel, R. **Learning the Stein discrepancy for training and evaluating energy-based models without sampling.** In ICML 2020.

**L1 IMQ & Cauchy RFF (Random Fourier Features):** Huggins, J. and Mackey, L. **Random feature Stein discrepancies.** In NeurIPS 2018.

# Thank you for your attention!

## Any Questions?



[MMDAagg](#)



[KSDAagg](#)



[AggInc](#)



[Code](#)