

# MMD Aggregated Two-Sample Test & KSD Aggregated Goodness-of-fit Test

MMDAgg: Antonin Schrab I. Kim M. Albert B. Laurent B. Guedj A. Gretton

KSDAgg: Antonin Schrab B. Guedj A. Gretton



UCL CENTRE FOR  
ARTIFICIAL INTELLIGENCE



Overview

## MMDAgg: Theoretical contributions

- aggregate MMD tests with different kernel bandwidths: no data splitting
- minimax adaptive over Sobolev balls: general kernels, estimated quantiles

## Two-sample problem

samples	$\mathbb{X}_m := (X_1, \dots, X_m)$	$X_i \stackrel{\text{iid}}{\sim} p$ in $\mathbb{R}^d$
samples	$\mathbb{Y}_n := (Y_1, \dots, Y_n)$	$Y_j \stackrel{\text{iid}}{\sim} q$ in $\mathbb{R}^d$
$\mathcal{H}_0: p = q$	against	$\mathcal{H}_a: p \neq q$

## Maximum Mean Discrepancy

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

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

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

## MMD test for fixed bandwidth $\lambda$

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

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

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

## MMDAgg for collection of bandwidths $\Lambda$

$\Delta_a^\Lambda(\mathbb{X}_m, \mathbb{Y}_n) := 1\left(\widehat{\text{MMD}}_\lambda^2(\mathbb{X}_m, \mathbb{Y}_n) > \widehat{q}_{1-u_\Lambda}^\lambda$  for some  $\lambda \in \Lambda\right)$   
with positive weights  $(w_\lambda)_{\lambda \in \Lambda}$  satisfying  $\sum_{\lambda \in \Lambda} w_\lambda \leq 1$  and correction

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

## Minimax adaptivity over Sobolev balls

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

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

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

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

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

Background

## KSDAgg: Theoretical contributions

- aggregate KSD tests with different kernel bandwidths: no data splitting
- uniform separation rate upper bound: general kernels, estimated quantiles

## Goodness-of-fit problem

model	density $p$	score $\nabla \log p(z)$
samples	$\mathbb{Z}_n := (Z_1, \dots, Z_n)$	$Z_i \stackrel{\text{iid}}{\sim} q$ in $\mathbb{R}^d$
$\mathcal{H}_0: p = q$	against	$\mathcal{H}_a: p \neq q$

## Kernel Stein Discrepancy

$$h_{p,\lambda}(x, y) := (\nabla \log p(x)^\top \nabla \log p(y)) k_\lambda(x, y) + \nabla \log p(y)^\top \nabla_1 k_\lambda(x, y) + \nabla \log p(x)^\top \nabla_2 k_\lambda(x, y) + \sum_{i=1}^d \frac{\partial}{\partial x_i \partial y_i} k_\lambda(x, y)$$

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

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

$$\widehat{\text{KSD}}_{p,\lambda}^2(\mathbb{Z}_n) := \frac{1}{n(n-1)} \sum_{1 \leq i \neq j \leq n} h_{p,\lambda}(Z_i, Z_j)$$

## KSD test for fixed bandwidth $\lambda$

$$\Delta_a^\lambda(\mathbb{Z}_n) := 1\left(\widehat{\text{KSD}}_{p,\lambda}^2(\mathbb{Z}_n) > \widehat{q}_{1-a}^\lambda\right)$$

**Quantile:**  $\widehat{q}_{1-a}^\lambda$  is  $[B(1-a)]$ -th largest of the  $B$   $\mathcal{H}_0$ -simulated test statistics

**Wild bootstrap:**  $\frac{1}{n(n-1)} \sum_{1 \leq i \neq j \leq n} \varepsilon_i \varepsilon_j h_{p,\lambda}(Z_i, Z_j)$   $\varepsilon_i \stackrel{\text{iid}}{\sim} \text{Unif}\{-1, 1\}$

**Parametric bootstrap:**  $\frac{1}{N(N-1)} \sum_{1 \leq i \neq j \leq N} h_{p,\lambda}(\tilde{Z}_i, \tilde{Z}_j)$   $\tilde{Z}_i \stackrel{\text{iid}}{\sim} p$

## KSDAgg for collection of bandwidths $\Lambda$

$\Delta_a^\Lambda(\mathbb{Z}_n) := 1\left(\widehat{\text{KSD}}_{p,\lambda}^2(\mathbb{Z}_n) > \widehat{q}_{1-u_\Lambda}^\lambda$  for some  $\lambda \in \Lambda\right)$   
with positive weights  $(w_\lambda)_{\lambda \in \Lambda}$  satisfying  $\sum_{\lambda \in \Lambda} w_\lambda \leq 1$  and correction

$$u_\Lambda = \sup\left\{u > 0 : \mathbb{P}_{p \times p}\left(\max_{\lambda \in \Lambda} \left(\widehat{\text{KSD}}_{p,\lambda}^2(\mathbb{Z}_n) - \widehat{q}_{1-uw_\lambda}^\lambda\right) > 0\right) \leq a\right\}$$

## Uniform separation rate

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

**Kernel assumption:**  $A_\lambda := \mathbb{E}_{q,q}[h_{p,\lambda}(Z, Z')^2] < \infty$

The condition

$$\|p - q\|_2^2 \geq \min_{\lambda \in \Lambda} \left(\|(p - q) - h_{p,\lambda} \diamond (p - q)\|_2^2 + C \ln\left(\frac{1}{aw_\lambda}\right) \frac{\sqrt{A_\lambda}}{\beta n}\right)$$

guarantees control over the probability of type II error of KSDAgg

$$\mathbb{P}_q(\Delta_{a,p}^\Lambda(\mathbb{Z}_n) = 0) \leq \beta.$$

Aggregation

## MMDAgg & KSDAgg: Experimental results

MMDAgg & KSDAgg obtain **higher power** than alternative state-of-the-art approaches to MMD & KSD kernel adaptation.

**Median bandwidth:**  $\lambda_{med}$  is the median  $L^2$ -distance between all the samples

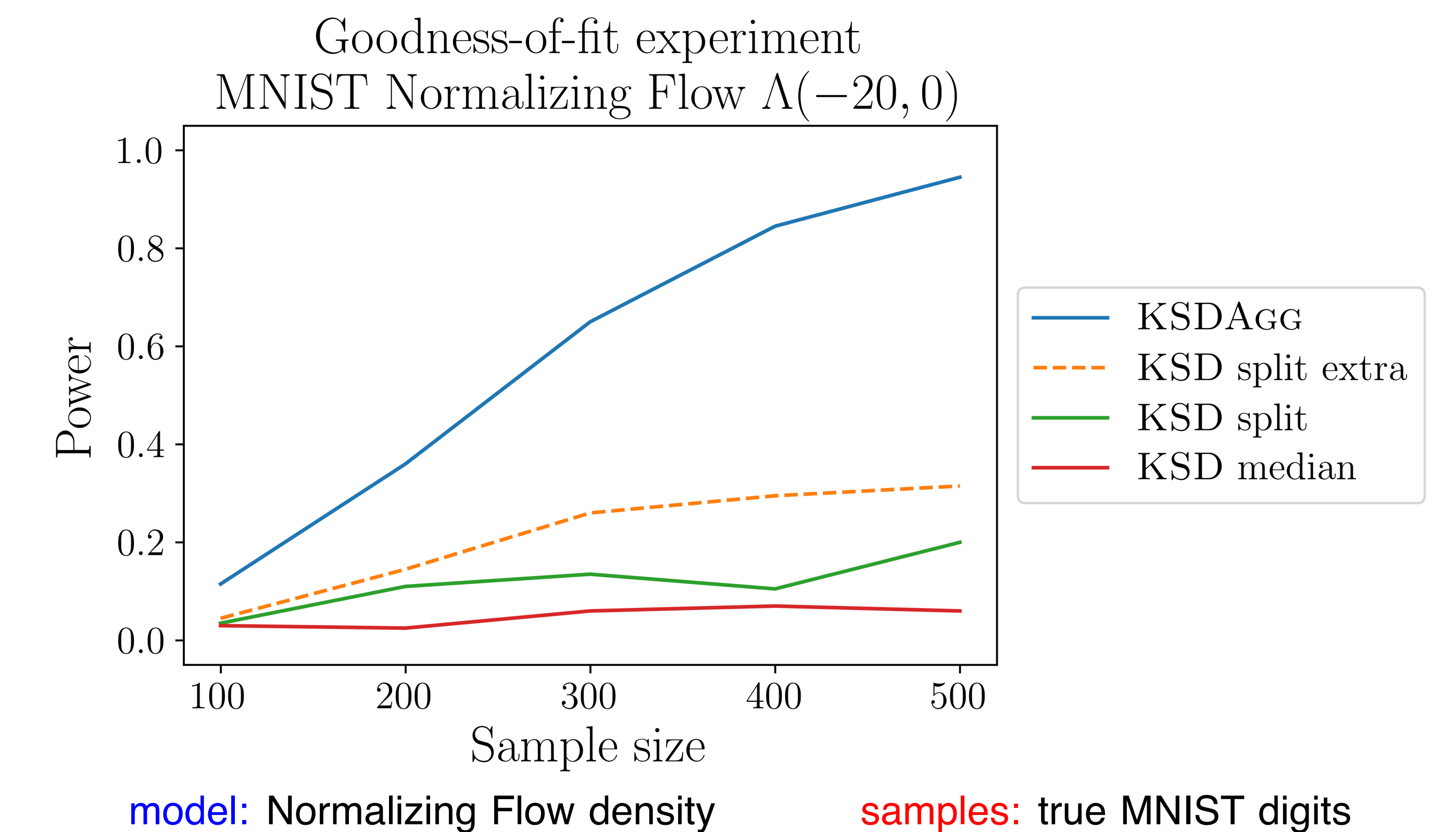
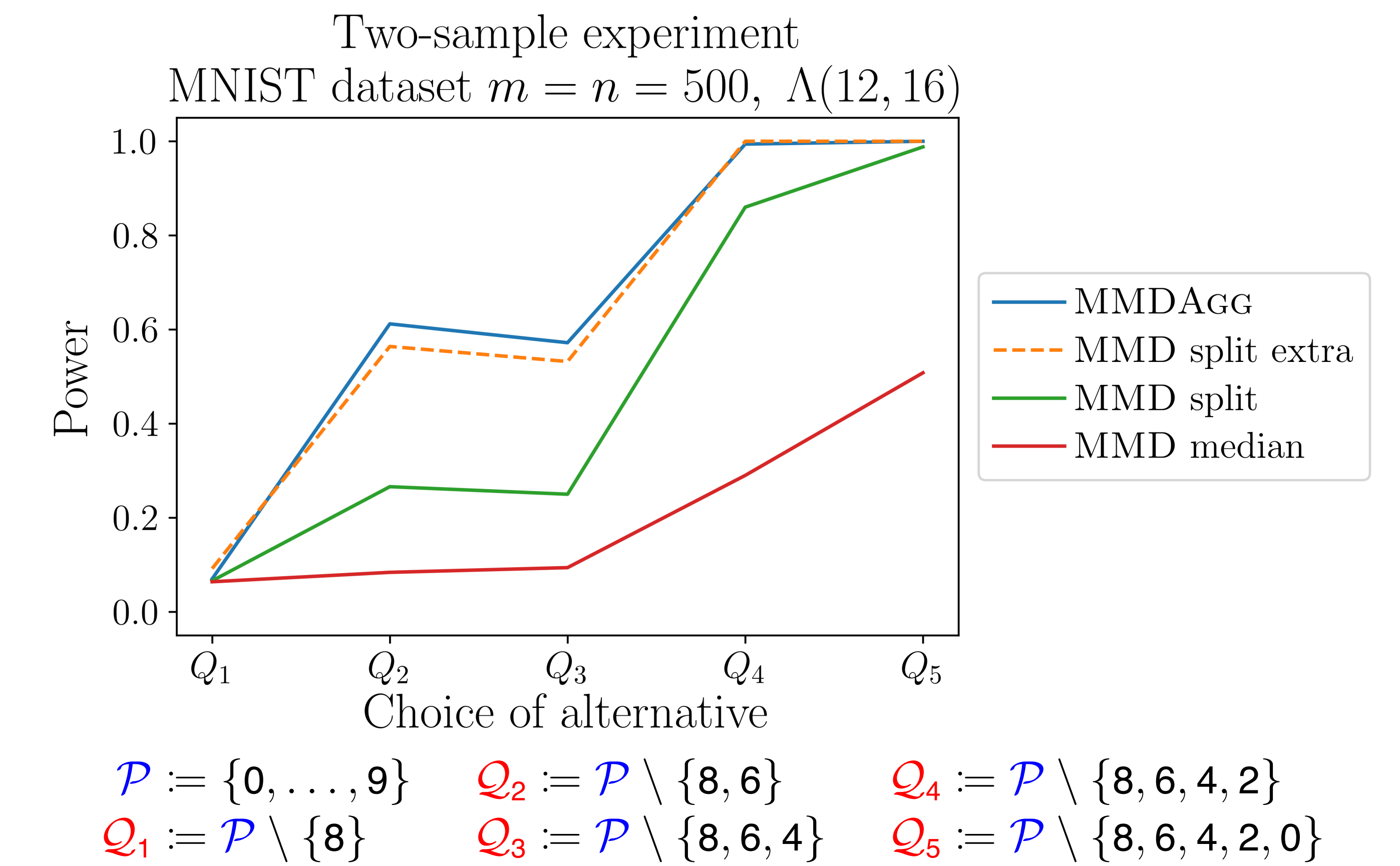
**Collection of bandwidths with uniform weights:**

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

**Split test:** split the data in two parts

- **1st part:** select bandwidth  $\lambda^*$  maximizing a proxy for asymptotic power
- **2nd part:** perform test with selected bandwidth  $\lambda^*$

**Split extra test:** select bandwidth  $\lambda^*$  using extra data



## MMDAgg



## KSDAgg



Theory