

Adaptive two-sample testing

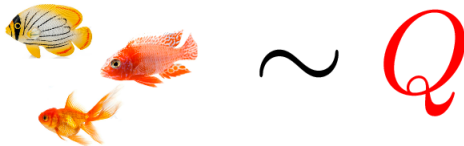
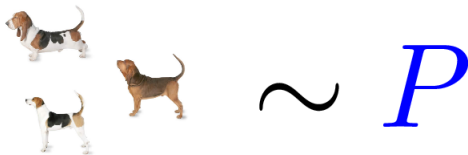
Arthur Gretton

Gatsby Computational Neuroscience Unit,
Google Deepmind

Cambridge, 2023

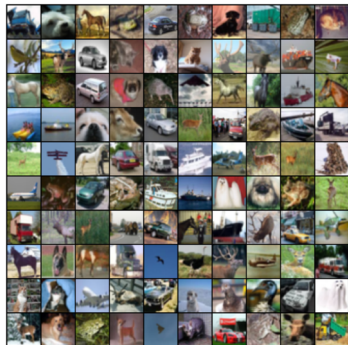
Comparing two samples

- Given: Samples from unknown distributions P and Q .
- Goal: do P and Q differ?

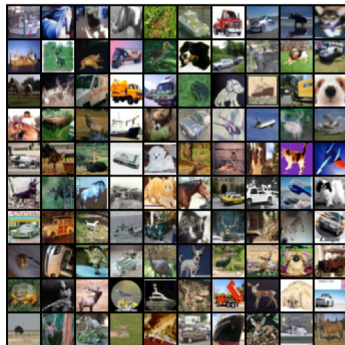


A real-life example: two-sample tests

- Goal: do P and Q differ?



CIFAR 10 samples



Cifar 10.1 samples

Significant difference?

Feng, Xu, Lu, Zhang, G., Sutherland, Learning Deep Kernels for Non-Parametric Two-Sample Tests, ICML 2020

Sutherland, Tung, Strathmann, De, Ramdas, Smola, G., ICLR 2017.

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

where $m \leq n$ and $n \leq Cm$.

Hypothesis test: function $\Delta_\alpha(\mathbb{X}_m, \mathbb{Y}_n)$

$\mathcal{H}_0: p = q$

against

$\mathcal{H}_1: p \neq q$

$\Delta_\alpha(\mathbb{X}_m, \mathbb{Y}_n) = 1$

\iff

reject \mathcal{H}_0

$\Delta_\alpha(\mathbb{X}_m, \mathbb{Y}_n) = 0$

\iff

fail to reject \mathcal{H}_0

Type II error β

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

Type I error: controlled by α by design

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

Outline

Two sample testing

- Test statistic: Maximum Mean Discrepancy (MMD)...
 - ...as a difference in feature means
 - ...as an integral probability metric
- Statistical testing with the MMD
- “How to choose the best kernel”
 - using aggregation (no sample splitting)
 - minimax guarantees with Sobolev smoothness assumption

Maximum Mean Discrepancy

Kernel methods, feature representation

Kernels: dot products of features

Feature map $\varphi(x) \in \mathcal{F}$,

$$\varphi(x) = [\dots \varphi_i(x) \dots] \in \ell_2$$

For positive definite k ,

$$k(x, x') = \langle \varphi(x), \varphi(x') \rangle_{\mathcal{F}}$$

Infinitely many features $\varphi(x)$, dot product in closed form!

Kernel methods, feature representation

Kernels: dot products of features

Feature map $\varphi(x) \in \mathcal{F}$,

$$\varphi(x) = [\dots \varphi_i(x) \dots] \in \ell_2$$

For positive definite k ,

$$k(x, x') = \langle \varphi(x), \varphi(x') \rangle_{\mathcal{F}}$$

Infinitely many features $\varphi(x)$, dot product in closed form!

Exponentiated quadratic kernel

$$k(x, x') = \exp \left(-\gamma \|x - x'\|^2 \right)$$

$$\varphi(x) = \begin{bmatrix} \varphi_1(x) \\ \varphi_2(x) \\ \varphi_3(x) \\ \vdots \end{bmatrix}$$

Features: Gaussian Processes for Machine learning, Rasmussen and Williams, Ch. 4.

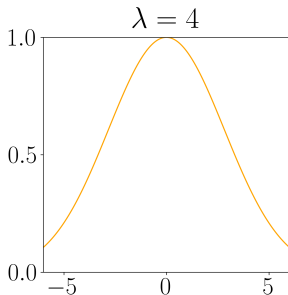
Kernels and bandwidths

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

Assumptions: K_1, \dots, K_d integrable (to 1) and square integrable

Eg: Gaussian ($K_i(u) \propto e^{-u^2}$), Laplace ($K_i(u) \propto e^{-|u|}$), Matérn, ...

Gaussian kernel: $k_{\lambda}(\mathbf{x} - \mathbf{y}) := \frac{1}{(\pi)^{d/2}} \prod_{i=1}^d \frac{1}{\lambda_i} \exp\left(-\frac{(\mathbf{x}_i - \mathbf{y}_i)^2}{\lambda_i^2}\right)$



Kernel mean embedding

Function evaluation in an RKHS:

$$f(\textcolor{red}{x}) = \langle f, \varphi_{\textcolor{red}{x}} \rangle_{\mathcal{F}}$$

Expectation evaluation in an RKHS:

$$\mathbb{E}_P(f(\textcolor{blue}{X})) = \mathbb{E}_P \langle f, \varphi_{\textcolor{blue}{X}} \rangle_{\mathcal{F}} = \langle f, \mathbb{E}_P \varphi_{\textcolor{blue}{X}} \rangle_{\mathcal{F}} =: \langle f, \mu_P \rangle_{\mathcal{F}}$$

as long as feature map Bochner integrable: $\mathbb{E}_P \|\varphi_{\textcolor{blue}{X}}\| = \mathbb{E}_P \sqrt{k_{\lambda}(\textcolor{blue}{X}, \textcolor{blue}{X})} < \infty$.

μ_P gives you **expectations** of all **RKHS functions**

“Kernel trick” for mean embeddings:

$$\langle \mu_{\textcolor{blue}{P}}, \mu_{\textcolor{red}{Q}} \rangle_{\mathcal{F}} = \mathbb{E}_{\textcolor{blue}{P}, \textcolor{red}{Q}} k_{\lambda}(\textcolor{blue}{X}, \textcolor{red}{Y})$$

for $\textcolor{blue}{X} \sim \textcolor{blue}{P}$ and $\textcolor{red}{Y} \sim \textcolor{red}{Q}$.

The maximum mean discrepancy

The maximum mean discrepancy is the distance between feature means:

$$\begin{aligned}\text{MMD}_{\lambda}^2(p, q) &= \|\mu_P - \mu_Q\|_{\mathcal{F}}^2 \\ &= \underbrace{\mathbb{E}_P k_{\lambda}(X, X')}_{(a)} + \underbrace{\mathbb{E}_Q k_{\lambda}(Y, Y')}_{(a)} - 2 \underbrace{\mathbb{E}_{P, Q} k_{\lambda}(X, Y)}_{(b)}\end{aligned}$$

(a)= within distrib. similarity, (b)= cross-distrib. similarity.

Characteristic kernels on \mathbb{R}^d

Characteristic kernel: $\text{MMD}_{\lambda}^2(p, q) = 0$ iff $p = q$ Fukumizu et al. [NIPS07b], Sriperumbudur et al. [COLT08]

When are **translation invariant** kernels $k_{\lambda}(x, y) = k_{\lambda}(x - y)$ characteristic on \mathbb{R}^d ?

Characteristic kernels on \mathbb{R}^d

Characteristic kernel: $\text{MMD}_{\lambda}^2(p, q) = 0$ iff $p = q$ Fukumizu et al. [NIPS07b], Sriperumbudur et al. [COLT08]

When are **translation invariant** kernels $k_{\lambda}(x, y) = k_{\lambda}(x - y)$ characteristic on \mathbb{R}^d ?

Bochner's theorem:

$$k_{\lambda}(x - y) = \int_{\mathbb{R}^d} e^{-i(x-y)^{\top} \omega} d\Lambda(\omega)$$

$\Lambda(\omega)$ finite non-negative Borel measure.

Characteristic function of P via **Fourier transform**

$$\varphi_P(\omega) = \int_{\mathbb{R}^d} e^{i x^{\top} \omega} dP(x)$$

Characteristic kernels on \mathbb{R}^d

Fourier representation of MMD on \mathbb{R}^d :

$$\text{MMD}_{\lambda}^2(\mathbf{p}, \mathbf{q}) = \int |\varphi_{\mathbf{P}}(\omega) - \varphi_{\mathbf{Q}}(\omega)|^2 d\Lambda(\omega)$$

Proof:

$$\begin{aligned} & \text{MMD}_{\lambda}^2(\mathbf{p}, \mathbf{q}) \\ &:= E_{\mathbf{P}} k_{\lambda}(x - x') + E_{\mathbf{Q}} k_{\lambda}(y - y') - 2E_{\mathbf{P}, \mathbf{Q}} k_{\lambda}(x, y) \\ &= \int \int \left[k_{\lambda}(s - t) d(\mathbf{P} - \mathbf{Q})(s) \right] d(\mathbf{P} - \mathbf{Q})(t) \\ &= \int_{\mathbb{R}^d} |\phi_{\mathbf{P}}(\omega) - \phi_{\mathbf{Q}}(\omega)|^2 d\Lambda(\omega) \end{aligned}$$

Characteristic kernels on \mathbb{R}^d

Fourier representation of MMD on \mathbb{R}^d :

$$\text{MMD}_{\lambda}^2(p, q) = \int |\varphi_P(\omega) - \varphi_Q(\omega)|^2 d\Lambda(\omega)$$

Proof:

$$\begin{aligned} & \text{MMD}_{\lambda}^2(p, q) \\ &:= E_P k_{\lambda}(x - x') + E_Q k_{\lambda}(y - y') - 2E_{P, Q} k_{\lambda}(x, y) \\ &= \int \int \left[k_{\lambda}(s - t) d(P - Q)(s) \right] d(P - Q)(t) \\ &= \int_{\mathbb{R}^d} |\phi_P(\omega) - \phi_Q(\omega)|^2 d\Lambda(\omega) \end{aligned}$$

Characteristic kernels on \mathbb{R}^d

Fourier representation of MMD on \mathbb{R}^d :

$$\text{MMD}_{\lambda}^2(\mathbf{p}, \mathbf{q}) = \int |\varphi_{\mathbf{P}}(\omega) - \varphi_{\mathbf{Q}}(\omega)|^2 d\Lambda(\omega)$$

Proof:

$$\begin{aligned} & \text{MMD}_{\lambda}^2(\mathbf{p}, \mathbf{q}) \\ &:= E_{\mathbf{P}} k_{\lambda}(x - x') + E_{\mathbf{Q}} k_{\lambda}(y - y') - 2E_{\mathbf{P}, \mathbf{Q}} k_{\lambda}(x, y) \\ &= \int \int \left[k_{\lambda}(s - t) d(\mathbf{P} - \mathbf{Q})(s) \right] d(\mathbf{P} - \mathbf{Q})(t) \\ &= \int_{\mathbb{R}^d} |\phi_{\mathbf{P}}(\omega) - \phi_{\mathbf{Q}}(\omega)|^2 d\Lambda(\omega) \end{aligned}$$

Summary: characteristic kernels on \mathbb{R}^d

A translation invariant k_λ is **characteristic** for prob. measures on \mathbb{R}^d if and only if

$$\text{supp}(\Lambda) = \mathbb{R}^d$$

(i.e. support zero on at most a countable set) Sriperumbudur et al. [COLT08, JMLR10]

Corollary: any continuous, compactly supported k_λ characteristic (since Fourier spectrum $\Lambda(\omega)$ cannot be zero on an interval).

1-D proof sketch from [Mallat, 99, Theorem 2.6], proof on \mathbb{R}^d via distribution theory in Sriperumbudur et al. [JMLR10, Corollary 10 p. 1535]

Summary: characteristic kernels on \mathbb{R}^d

A translation invariant k_λ is **characteristic** for prob. measures on \mathbb{R}^d if and only if

$$\text{supp}(\Lambda) = \mathbb{R}^d$$

(i.e. support zero on at most a countable set) Sriperumbudur et al. [COLT08, JMLR10]

Corollary: any continuous, compactly supported k_λ characteristic (since Fourier spectrum $\Lambda(\omega)$ cannot be zero on an interval).

1-D proof sketch from [Mallat, 99, Theorem 2.6], proof on \mathbb{R}^d via distribution theory in Sriperumbudur et al. [JMLR10, Corollary 10 p. 1535]

Two-Sample Testing with MMD

A statistical test using MMD

The empirical MMD:

$$\widehat{\text{MMD}}_{\lambda}^2(\mathbb{X}_m, \mathbb{Y}_n) = \frac{1}{m(m-1)} \sum_{i \neq j} k_{\lambda}(x_i, x_j) + \frac{1}{n(n-1)} \sum_{i \neq j} k_{\lambda}(y_i, y_j) \\ - \frac{2}{mn} \sum_{i,j} k_{\lambda}(x_i, y_j)$$

Two-sample MMD test:

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

A statistical test using MMD

The empirical MMD:

$$\widehat{\text{MMD}}_{\lambda}^2(\mathbb{X}_m, \mathbb{Y}_n) = \frac{1}{m(m-1)} \sum_{i \neq j} k_{\lambda}(\mathbf{x}_i, \mathbf{x}_j) + \frac{1}{n(n-1)} \sum_{i \neq j} k_{\lambda}(\mathbf{y}_i, \mathbf{y}_j) \\ - \frac{2}{mn} \sum_{i,j} k_{\lambda}(\mathbf{x}_i, \mathbf{y}_j)$$

Two-sample MMD test:

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

Want threshold $q_{1-\alpha}^{\lambda}$ for test $\widehat{\text{MMD}}_{\lambda}^2(\mathbb{X}_m, \mathbb{Y}_n)$ to get false positive rate α

Asymptotics of \widehat{MMD}^2 when $P = Q$

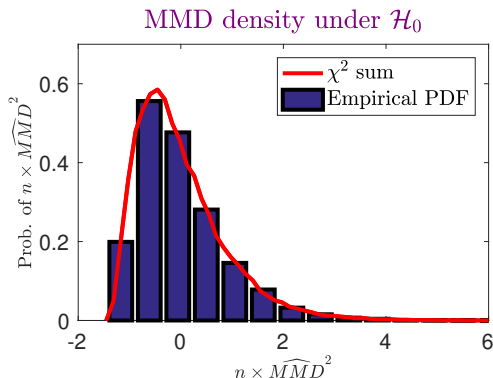
$P = Q = \mathcal{N}(0, 1)$, statistic has asymptotic distribution

$$(m+n)\widehat{MMD}_{\lambda}^2(\mathbb{X}_m, \mathbb{Y}_n) \sim \sum_{l=1}^{\infty} \lambda_l [z_l^2 - 2]$$

where

$$\lambda_i \psi_i(x') = \int_{\mathcal{X}} \underbrace{\tilde{k}_{\lambda}(x, x')}_{\text{centred}} \psi_i(x) dP(x)$$

$$z_l \sim \mathcal{N}(0, 2) \quad \text{i.i.d.}$$



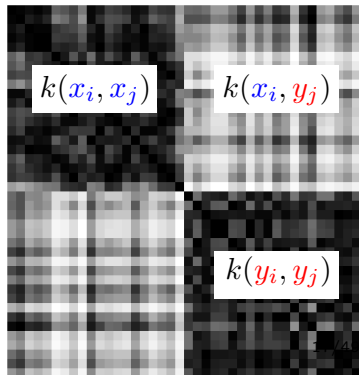
How do we get the test threshold $q_{1-\alpha}^\lambda$?

Original empirical MMD for dogs and fish:

$$X = \left[\text{dog1} \quad \text{dog2} \quad \text{dog3} \quad \dots \right]$$

$$Y = \left[\text{fish1} \quad \text{fish2} \quad \text{fish3} \quad \dots \right]$$

$$\begin{aligned} \widehat{MMD}^2 &= \frac{1}{n(n-1)} \sum_{i \neq j} k(\mathbf{x}_i, \mathbf{x}_j) \\ &\quad + \frac{1}{n(n-1)} \sum_{i \neq j} k(\mathbf{y}_i, \mathbf{y}_j) \\ &\quad - \frac{2}{n^2} \sum_{i,j} k(\mathbf{x}_i, \mathbf{y}_j) \end{aligned}$$



How do we get test threshold $q_{1-\alpha}^\lambda$?

Permuted dog and fish samples (merdogs):

$$\tilde{X} = \left[\text{fish} \quad \text{dog} \quad \text{fish} \quad \dots \right]$$

$$\tilde{Y} = \left[\text{dog} \quad \text{fish} \quad \text{dog} \quad \dots \right]$$



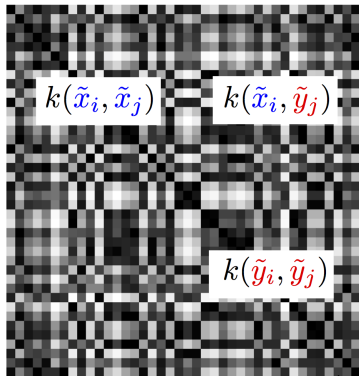
How do we get test threshold $q_{1-\alpha}^\lambda$?

Permuted dog and fish samples (merdogs):

$$\tilde{X} = \left[\text{fish} \quad \text{dog} \quad \text{fish} \quad \dots \right]$$

$$\tilde{Y} = \left[\text{dog} \quad \text{fish} \quad \text{dog} \quad \dots \right]$$

$$\begin{aligned} & \widehat{\text{MMD}}_\lambda^2(\mathbb{X}_m, \mathbb{Y}_n) \\ &= \frac{1}{m(m-1)} \sum_{i \neq j} k(\tilde{x}_i, \tilde{x}_j) \\ &+ \frac{1}{n(n-1)} \sum_{i \neq j} k(\tilde{y}_i, \tilde{y}_j) \\ &- \frac{2}{mn} \sum_{i,j} k(\tilde{x}_i, \tilde{y}_j) \end{aligned}$$



MMD test thresholds: permutation, wild bootstrap

Two-sample MMD test:

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

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 $(\mathbb{X}_m^{\sigma}, \mathbb{Y}_n^{\sigma}) = \sigma(\mathbb{X}_m \cup \mathbb{Y}_n)$

MMD test thresholds: permutation, wild bootstrap

Two-sample MMD test:

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

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 $(\mathbb{X}_m^{\sigma}, \mathbb{Y}_n^{\sigma}) = \sigma(\mathbb{X}_m \cup \mathbb{Y}_n)$

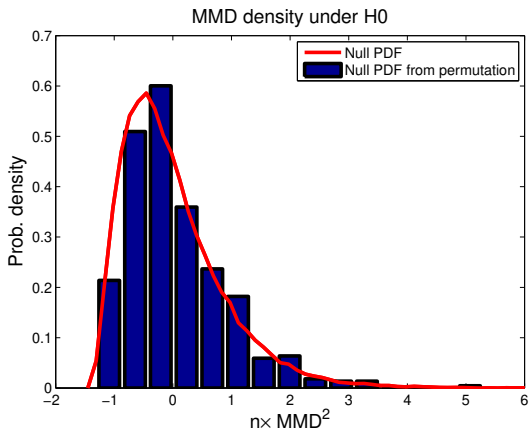
Non-asymptotic level (permutation): $\mathbb{P}_{p \times p}(\Delta_{\alpha}^{\lambda, B}(\mathbb{X}_m, \mathbb{Y}_n) = 1) \leq \alpha$,

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

Approx. null distribution of \widehat{MMD}^2 via permutation

Null distribution estimated from 500 permutations

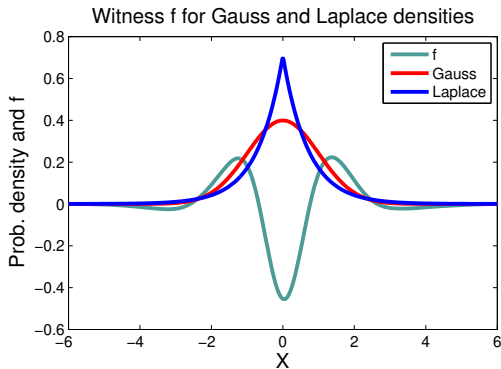
Example: $P = Q = \mathcal{N}(0, 1)$



Kernel choice: MMD as an IPM

Maximum mean discrepancy: smooth function for P vs Q

$$\text{MMD}_\lambda(p, q) := \sup_{\|f\|_{\mathcal{F}} \leq 1} [\mathbb{E}_P f(X) - \mathbb{E}_Q f(Y)]$$



Kernel Choice: MMD as an IPM

- Simple choice: Gaussian

$$k_{\lambda}(x, y) = \frac{1}{(\pi)^{d/2}} \prod_{i=1}^d \frac{1}{\lambda_i} \exp\left(-\frac{(x_i - y_i)^2}{\lambda_i^2}\right)$$

- *Characteristic*: for any σ : for any P and Q , power $\rightarrow 1$ as $n \rightarrow \infty$

Kernel Choice: MMD as an IPM

- Simple choice: Gaussian

$$k_{\lambda}(x, y) = \frac{1}{(\pi)^{d/2}} \prod_{i=1}^d \frac{1}{\lambda_i} \exp\left(-\frac{(x_i - y_i)^2}{\lambda_i^2}\right)$$

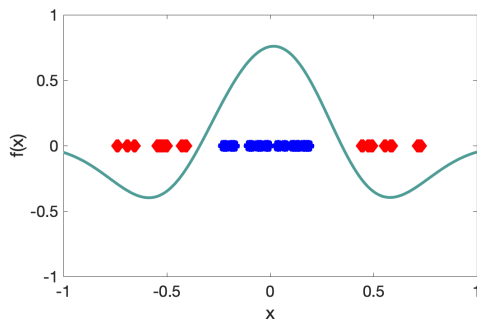
- *Characteristic*: for any σ : for any P and Q , power $\rightarrow 1$ as $n \rightarrow \infty$
- But choice of $\lambda_1 \cdots \lambda_d$ is very important for finite $m, n \dots$

Kernel Choice: MMD as an IPM

- Simple choice: Gaussian

$$k_{\lambda}(x, y) = \frac{1}{(\pi)^{d/2}} \prod_{i=1}^d \frac{1}{\lambda_i} \exp\left(-\frac{(x_i - y_i)^2}{\lambda_i^2}\right)$$

- *Characteristic*: for any σ : for any P and Q , power $\rightarrow 1$ as $n \rightarrow \infty$
- But choice of $\lambda_1 \cdots \lambda_d$ is very important for finite $m, n \dots$

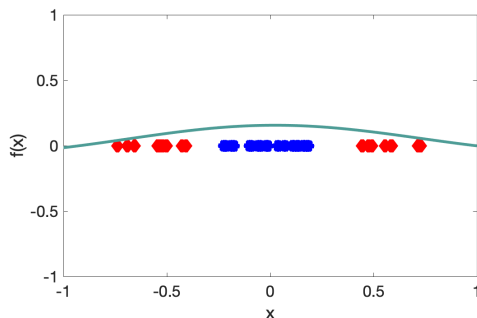


Kernel Choice: MMD as an IPM

- Simple choice: Gaussian

$$k_{\lambda}(x, y) = \frac{1}{(\pi)^{d/2}} \prod_{i=1}^d \frac{1}{\lambda_i} \exp\left(-\frac{(x_i - y_i)^2}{\lambda_i^2}\right)$$

- *Characteristic:* for any σ : for any P and Q , power $\rightarrow 1$ as $n \rightarrow \infty$
- But choice of $\lambda_1 \cdots \lambda_d$ is very important for finite $m, n \dots$

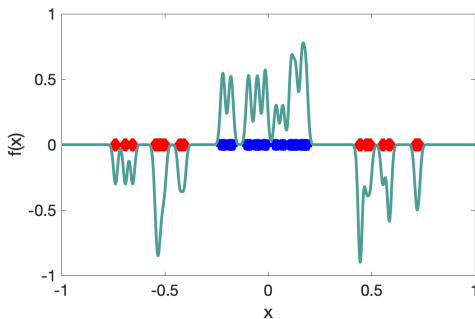


Kernel Choice: MMD as an IPM

- Simple choice: Gaussian

$$k_{\lambda}(x, y) = \frac{1}{(\pi)^{d/2}} \prod_{i=1}^d \frac{1}{\lambda_i} \exp\left(-\frac{(x_i - y_i)^2}{\lambda_i^2}\right)$$

- *Characteristic:* for any σ : for any P and Q , power $\rightarrow 1$ as $n \rightarrow \infty$
- But choice of $\lambda_1 \cdots \lambda_d$ is very important for finite $m, n \dots$



Test power for known
smoothness of $p - q$

Sobolev balls

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

Sobolev balls:

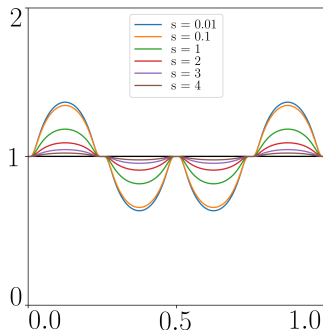
$$\mathcal{S}_d^s(R) := \left\{ f \in L^1(\mathbb{R}^d) \cap L^2(\mathbb{R}^d) : \int_{\mathbb{R}^d} \|\xi\|_2^{2s} |\hat{f}(\xi)|^2 d\xi \leq (2\pi)^d R^2 \right\}$$

radius $R > 0$

dimension d

smoothness parameter $s > 0$

Fourier transform $\hat{f}(\xi) := \int_{\mathbb{R}^d} f(x) e^{-ix^\top \xi} dx$



Sobolev balls

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

Sobolev balls:

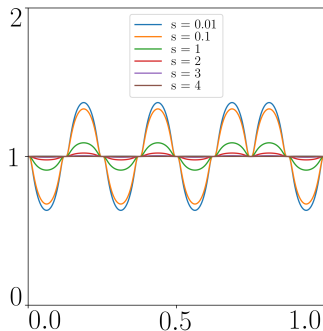
$$\mathcal{S}_d^s(R) := \left\{ f \in L^1(\mathbb{R}^d) \cap L^2(\mathbb{R}^d) : \int_{\mathbb{R}^d} \|\xi\|_2^{2s} |\hat{f}(\xi)|^2 d\xi \leq (2\pi)^d R^2 \right\}$$

radius $R > 0$

dimension d

smoothness parameter $s > 0$

Fourier transform $\hat{f}(\xi) := \int_{\mathbb{R}^d} f(x) e^{-ix^\top \xi} dx$



MMD test power, known smoothness

Theorem (MMD test minimax optimality)

For *known* smoothness s , assuming $p - q \in S_d^s(R)$ and setting

$$\lambda_i^* := (m + n)^{-2/(4s+d)}$$

for $i = 1, \dots, d$, the condition

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

guarantees control of the type II error of the MMD test

$$\mathbb{P}_{p \times q \times r} \left(\Delta_{\alpha}^{\lambda, B}(\mathbb{X}_m, \mathbb{Y}_n) = 0 \right) \leq \beta.$$

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

Proof of theorem 1 (next few slides)

For $\alpha, \beta \in (0, 1)$, Type II error control

$$\mathbb{P}_{\mathbf{p} \times \mathbf{q} \times \mathbf{r}} \left(\Delta_{\alpha}^{\lambda, B}(\mathbb{X}_m, \mathbb{Y}_n) = 0 \right) \leq \beta$$

is implied by (Chebyshev)

$$\mathbb{P}_{\mathbf{p} \times \mathbf{q} \times \mathbf{r}} \left(\text{MMD}_{\lambda}^2(\mathbf{p}, \mathbf{q}) \geq \sqrt{\frac{2}{\beta} \text{var}_{\mathbf{p} \times \mathbf{q}} \left(\widehat{\text{MMD}}_{\lambda}^2(\mathbb{X}_m, \mathbb{Y}_n) \right)} + \widehat{q}_{1-\alpha}^{\lambda} \right) \geq 1 - \frac{\beta}{2}$$

Proof of theorem 1 (next few slides)

$$\underbrace{\text{MMD}_{\lambda}^2(\mathbf{p}, \mathbf{q})}_{(A)} \geq \underbrace{\sqrt{\frac{2}{\beta} \text{var}_{\mathbf{p} \times \mathbf{q}} \left(\widehat{\text{MMD}}_{\lambda}^2(\mathbf{X}_m, \mathbf{Y}_n) \right)}}_{(B)} + \underbrace{\widehat{q}_{1-\alpha}^{\lambda}}_{(C)}$$

We address each of the three terms (A), (B), (C) in turn.

Breakdown of the MMD (A)

The MMD can be decomposed

$$\begin{aligned}
 \text{MMD}_{\lambda}^2(p, q) &= \langle p - q, k_{\lambda} * (p - q) \rangle_2 \\
 &= \frac{1}{2} \left(\|p - q\|_2^2 + \|k_{\lambda} * (p - q)\|_2^2 \right. \\
 &\quad \left. - \|k_{\lambda} * (p - q) - (p - q)\|_2^2 \right)
 \end{aligned}$$

Breakdown of the MMD (A)

The MMD can be decomposed

$$\begin{aligned}\text{MMD}_{\lambda}^2(p, q) &= \langle p - q, k_{\lambda} * (p - q) \rangle_2 \\ &= \frac{1}{2} \left(\|p - q\|_2^2 + \|k_{\lambda} * (p - q)\|_2^2 \right. \\ &\quad \left. - \|k_{\lambda} * (p - q) - (p - q)\|_2^2 \right)\end{aligned}$$

- Keep the first term (test “radius” for power $1 - \beta$): $\|p - q\|_2^2$
- Get rid of second term using variance (next slides): $\|k_{\lambda} * (p - q)\|_2^2$
- Bound the final term: if $p - q \in \mathcal{S}_d^s(R)$, then $\exists S \in (0, 1)$ such that

$$\|k_{\lambda} * (p - q) - (p - q)\|_2^2 - S^2 \|p - q\|_2^2 \leq C_0(d, s, R) \sum_{i=1}^d \lambda_i^{2s}$$

Updating the power condition after (A)

The power condition (which needs to hold with probability $1 - \beta/2$)

$$\underbrace{\text{MMD}_{\lambda}^2(p, q)}_{(A)} \geq \underbrace{\sqrt{\frac{2}{\beta} \text{var}_{p \times q} \left(\widehat{\text{MMD}}_{\lambda}^2(\mathbb{X}_m, \mathbb{Y}_n) \right)}}_{(B)} + \underbrace{\widehat{q}_{1-\alpha}^{\lambda}}_{(C)},$$

Updating the power condition after (A)

The power condition (which needs to hold with probability $1 - \beta/2$)

$$\underbrace{\text{MMD}_{\lambda}^2(\mathbf{p}, \mathbf{q})}_{(A)} \geq \underbrace{\sqrt{\frac{2}{\beta} \text{var}_{\mathbf{p} \times \mathbf{q}} \left(\widehat{\text{MMD}}_{\lambda}^2(\mathbb{X}_m, \mathbb{Y}_n) \right)}}_{(B)} + \underbrace{\widehat{q}_{1-\alpha}^{\lambda}}_{(C)},$$

after updating (A), then becomes

$$\begin{aligned} (1 - S^2) \|\mathbf{p} - \mathbf{q}\|_2^2 &\geq C_0 \sum_{i=1}^d \lambda_i^{2s} - \|k_{\lambda} * (\mathbf{p} - \mathbf{q})\|_2^2 \\ &\quad + 2 \underbrace{\sqrt{\frac{2}{\beta} \text{var}_{\mathbf{p} \times \mathbf{q}} \left(\widehat{\text{MMD}}_{\lambda}^2(\mathbb{X}_m, \mathbb{Y}_n) \right)}}_{(B)} + 2 \underbrace{\widehat{q}_{1-\alpha}^{\lambda}}_{(C)} \end{aligned}$$

Bound on the variance (B)

Assume that $\max(\|\mathbf{p}\|_\infty, \|\mathbf{q}\|_\infty) \leq M$ for some $M > 0$.

$$\begin{aligned} & \text{var}_{\mathbf{p} \times \mathbf{q}} \left(\widehat{\text{MMD}}_{\lambda}^2(\mathbb{X}_m, \mathbb{Y}_n) \right) \\ & \leq C_1(M, d) \left(\frac{\|k_{\lambda} * (\mathbf{p} - \mathbf{q})\|_2^2}{m + n} + \frac{1}{(m + n)^2 \lambda_1 \cdots \lambda_d} \right). \end{aligned}$$

Bound on the variance (B)

Assume that $\max(\|p\|_\infty, \|q\|_\infty) \leq M$ for some $M > 0$.

$$\begin{aligned} & \text{var}_{p \times q} \left(\widehat{\text{MMD}}_\lambda^2(\mathbb{X}_m, \mathbb{Y}_n) \right) \\ & \leq C_1(M, d) \left(\frac{\|k_\lambda * (p - q)\|_2^2}{m + n} + \frac{1}{(m + n)^2 \lambda_1 \cdots \lambda_d} \right). \end{aligned}$$

Assuming $\lambda_1 \cdots \lambda_d \leq 1$,

$$\begin{aligned} (B) &= 2 \sqrt{\frac{2}{\beta} \text{var}_{p \times q} \left(\widehat{\text{MMD}}_\lambda^2(\mathbb{X}_m, \mathbb{Y}_n) \right)} \\ &\leq \|k_\lambda * (p - q)\|_2^2 + \frac{C'_1}{\sqrt{\beta}(m + n)\sqrt{\lambda_1 \cdots \lambda_d}}. \end{aligned}$$

Term $\|k_\lambda * (p - q)\|_2^2$ will cancel in the power condition.

Updating the power condition after (A), (B)

The power condition (which needs to hold with probability $1 - \beta/2$)

$$(1 - S^2) \|\mathbf{p} - \mathbf{q}\|_2^2 \geq C_0 \sum_{i=1}^d \lambda_i^{2s} - \|\mathbf{k}_\lambda * (\mathbf{p} - \mathbf{q})\|_2^2 \\ + 2 \underbrace{\sqrt{\frac{2}{\beta} \text{var}_{\mathbf{p} \times \mathbf{q}} \left(\widehat{\text{MMD}}_{\lambda}^2(\mathbb{X}_m, \mathbb{Y}_n) \right)}}_{(B)} + 2 \underbrace{\widehat{q}_{1-\alpha}^\lambda}_{(C)}.$$

Updating the power condition after (A), (B)

The power condition (which needs to hold with probability $1 - \beta/2$)

$$(1 - S^2) \| \textcolor{blue}{p} - \textcolor{red}{q} \|_2^2 \geq C_0 \sum_{i=1}^d \lambda_i^{2s} - \| k_\lambda * (\textcolor{blue}{p} - \textcolor{red}{q}) \|_2^2 \\ + 2 \underbrace{\sqrt{\frac{2}{\beta} \text{var}_{\textcolor{blue}{p} \times \textcolor{red}{q}} \left(\widehat{\text{MMD}}_{\lambda}^2(\textcolor{blue}{X}_m, \textcolor{red}{Y}_n) \right)}}_{(B)} + 2 \underbrace{\widehat{q}_{1-\alpha}^\lambda}_{(C)}.$$

After bounding (B), becomes

$$(1 - S^2) \| \textcolor{blue}{p} - \textcolor{red}{q} \|_2^2 \geq \frac{C'_1}{\sqrt{\beta}(\textcolor{blue}{m} + \textcolor{red}{n}) \lambda_1 \cdots \lambda_d} + C_0 \sum_{i=1}^d \lambda_i^{2s} + 2 \underbrace{\widehat{q}_{1-\alpha}^\lambda}_{(C)}.$$

Bound on estimated $1 - \alpha$ quantile (C)

Assume that $\max(\|\mathbf{p}\|_\infty, \|\mathbf{q}\|_\infty) \leq M$ for some $M > 0$. We have

$$\mathbb{P}_{\mathbf{p} \times \mathbf{q} \times \mathbf{r}} \left(\widehat{q}_{1-\alpha}^{\lambda} \leq C_2(M, d) \frac{\ln\left(\frac{1}{\alpha}\right)}{\sqrt{\beta}(\mathbf{m} + \mathbf{n})\sqrt{\lambda_1 \cdots \lambda_d}} \right) \geq 1 - \frac{\beta}{2}$$

for $B \geq \frac{3}{\alpha^2} \left(\ln\left(\frac{8}{\beta}\right) + \alpha(1 - \alpha) \right)$ and $\alpha \in (0, 0.5)$.

Updating the power condition after (A), (B), (C)

The power condition (which needs to hold with probability $1 - \beta/2$)

$$(1 - S^2) \| \textcolor{blue}{p} - \textcolor{red}{q} \|_2^2 \geq \frac{C'_1}{\sqrt{\beta}(\textcolor{blue}{m} + \textcolor{red}{n})\sqrt{\lambda_1 \cdots \lambda_d}} + C_0 \sum_{i=1}^d \lambda_i^{2s} + 2 \underbrace{\widehat{q}_{1-\alpha}^\lambda}_{(C)}.$$

Fine print: $\alpha \in (0, e^{-1})$, $B \geq \frac{3}{\alpha^2} \left(\ln \left(\frac{8}{\beta} \right) + \alpha(1 - \alpha) \right)$, and $\lambda_1 \cdots \lambda_d \leq 1$.

Updating the power condition after (A), (B), (C)

The power condition (which needs to hold with probability $1 - \beta/2$)

$$(1 - S^2) \|p - q\|_2^2 \geq \frac{C'_1}{\sqrt{\beta}(m + n)\sqrt{\lambda_1 \cdots \lambda_d}} + C_0 \sum_{i=1}^d \lambda_i^{2s} + 2 \underbrace{\widehat{q}_{1-\alpha}^\lambda}_{(C)}.$$

After updating (C)

$$\|p - q\|_2^2 \geq \frac{C_4(M, d, s, S, R)}{\sqrt{\beta}} \left(\sum_{i=1}^d \lambda_i^{2s} + \frac{\ln\left(\frac{1}{\alpha}\right)}{(m + n)\sqrt{\lambda_1 \cdots \lambda_d}} \right).$$

Fine print: $\alpha \in (0, e^{-1})$, $B \geq \frac{3}{\alpha^2} \left(\ln\left(\frac{8}{\beta}\right) + \alpha(1 - \alpha) \right)$, and $\lambda_1 \cdots \lambda_d \leq 1$.

Updating the power condition after (A), (B), (C)

The power condition (which needs to hold with probability $1 - \beta/2$)

$$(1 - S^2) \|p - q\|_2^2 \geq \frac{C'_1}{\sqrt{\beta}(\mathbf{m} + \mathbf{n})\sqrt{\lambda_1 \cdots \lambda_d}} + C_0 \sum_{i=1}^d \lambda_i^{2s} + 2 \underbrace{\widehat{q}_{1-\alpha}^\lambda}_{(C)}.$$

After updating (C)

$$\|p - q\|_2^2 \geq \frac{C_4(M, d, s, S, R)}{\sqrt{\beta}} \left(\sum_{i=1}^d \lambda_i^{2s} + \frac{\ln\left(\frac{1}{\alpha}\right)}{(\mathbf{m} + \mathbf{n})\sqrt{\lambda_1 \cdots \lambda_d}} \right).$$

Picking $\lambda_i^* := (\mathbf{m} + \mathbf{n})^{-2/(4s+d)}$ controls the Type II error when

$$\|p - q\|_2 \geq \frac{C}{\sqrt{\beta}} (\mathbf{m} + \mathbf{n})^{-2s/(4s+d)}$$

Fine print: $\alpha \in (0, e^{-1})$, $B \geq \frac{3}{\alpha^2} \left(\ln\left(\frac{8}{\beta}\right) + \alpha(1 - \alpha) \right)$, and $\lambda_1 \cdots \lambda_d \leq 1$.

Optimizing kernel parameters: aggregation

MMDAgg for a collection of bandwidths Λ

MMDAgg (MMD Aggregation): non-asymptotic level α

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

positive weights $(w_{\lambda})_{\lambda \in \Lambda}$ satisfying $\sum_{\lambda \in \Lambda} w_{\lambda} \leq 1$

Correction u_{α} 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 $\mathcal{O}(|\Lambda| (B_1 + B_2) (m + n)^2)$

MMDAgg for a collection of bandwidths Λ

MMDAgg (MMD Aggregation): non-asymptotic level α

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

positive weights $(w_{\lambda})_{\lambda \in \Lambda}$ satisfying $\sum_{\lambda \in \Lambda} w_{\lambda} \leq 1$

Correction u_{α} defined as

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

more powerful than Bonferroni correction as $u_{\alpha} \geq \alpha$

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

MMDAgg for a collection of bandwidths Λ

MMDAgg (MMD Aggregation): non-asymptotic level α

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

positive weights $(w_{\lambda})_{\lambda \in \Lambda}$ satisfying $\sum_{\lambda \in \Lambda} w_{\lambda} \leq 1$

Correction u_{α} 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 $\mathcal{O}(|\Lambda| (B_1 + B_2) (m + n)^2)$

MMDAgg for a collection of bandwidths Λ

MMDAgg (MMD Aggregation): non-asymptotic level α

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

positive weights $(w_{\lambda})_{\lambda \in \Lambda}$ satisfying $\sum_{\lambda \in \Lambda} w_{\lambda} \leq 1$

Correction u_{α} defined as

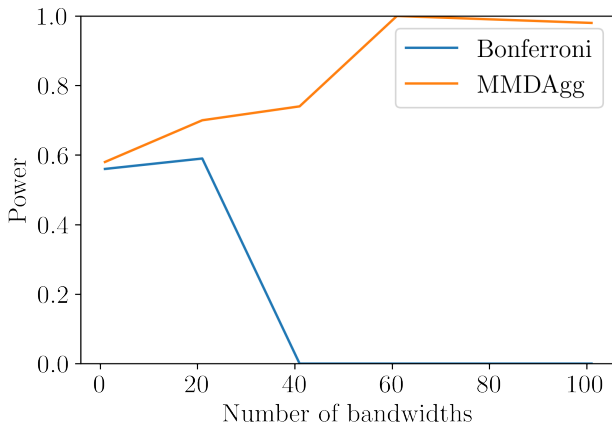
$$\sup \left\{ u > 0 : \mathbb{P}_{p \times p} \left(\max_{\lambda \in \Lambda} \left(\widehat{\text{MMD}}_{\lambda}^2(\mathbb{X}_m, \mathbb{Y}_n) - \hat{q}_{1-uw_{\lambda}}^{\lambda} \right) > 0 \right) \leq \alpha \right\}$$

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

Simple example: 3-d Gaussians with different means



$$\Lambda(i) := \left\{ 2^\ell \lambda_{\text{med}} : \ell \in \{-i, \dots, i\} \right\} \text{ for } i \in \{0, 10, 20, 30, 40, 50\}$$

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

MMDAgg test power guarantee

Theorem (MMDAgg minimax adaptivity)

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

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

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

guarantees control of the type II error of MMDAgg

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

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

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

Unlike the MMD test $\Delta_\alpha^{\lambda^*}$, MMDAgg $\Delta_\alpha^{\Lambda^*}$ is independent of s

MMDAgg parameter-free user-friendly implementation

Radial basis function (RBF) kernel: $k_{\lambda}(\mathbf{x}, \mathbf{y}) := K\left(\left\|\frac{\mathbf{x} - \mathbf{y}}{\lambda}\right\|\right)$

Collection of bandwidths Λ : discretisation of the interval $[\lambda_{\min}, \lambda_{\max}]$ where λ_{\min} and λ_{\max} are the (robust) minimum and maximum of $\left\{\|\mathbf{x} - \mathbf{y}\| : \mathbf{x} \in \mathbb{X}_m, \mathbf{y} \in \mathbb{Y}_n\right\}$

Possible to aggregate several kernels each with multiple bandwidths

Uniform weights: $w_{\lambda} := 1 / |\Lambda|$

Number of permutations / wild bootstraps: $B_1 = B_2 = 2000$

JAX: runs on either CPU or GPU (significant speed improvements)

■ JAX GPU runs 100 times faster than Numpy CPU

mmdagg package: github.com/antoninschrab/mmdagg

```
from mmdagg import mmdagg                # X shape (m, d)
output = mmdagg(X, Y) # 0 or 1            # Y shape (n, d)
```


MMDAgg parameter-free user-friendly implementation

Radial basis function (RBF) kernel: $k_{\lambda}(\mathbf{x}, \mathbf{y}) := K\left(\left\|\frac{\mathbf{x} - \mathbf{y}}{\lambda}\right\|\right)$

Collection of bandwidths Λ : discretisation of the interval $[\lambda_{\min}, \lambda_{\max}]$ where λ_{\min} and λ_{\max} are the (robust) minimum and maximum of $\left\{\|\mathbf{x} - \mathbf{y}\| : \mathbf{x} \in \mathbb{X}_m, \mathbf{y} \in \mathbb{Y}_n\right\}$

Possible to aggregate several kernels each with multiple bandwidths

Uniform weights: $w_{\lambda} := 1 / |\Lambda|$

Number of permutations / wild bootstraps: $B_1 = B_2 = 2000$

JAX: runs on either CPU or GPU (significant speed improvements)

■ JAX GPU runs 100 times faster than Numpy CPU

mmdagg package: github.com/antoninschrab/mmdagg

```
from mmdagg import mmdagg          # X shape (m, d)
output = mmdagg(X, Y) # 0 or 1      # Y shape (n, d)
```

MMDAgg parameter-free user-friendly implementation

Radial basis function (RBF) kernel: $k_{\lambda}(\mathbf{x}, \mathbf{y}) := K\left(\left\|\frac{\mathbf{x} - \mathbf{y}}{\lambda}\right\|\right)$

Collection of bandwidths Λ : discretisation of the interval $[\lambda_{\min}, \lambda_{\max}]$ where λ_{\min} and λ_{\max} are the (robust) minimum and maximum of $\{\|\mathbf{x} - \mathbf{y}\| : \mathbf{x} \in \mathbb{X}_m, \mathbf{y} \in \mathbb{Y}_n\}$

Possible to aggregate several kernels each with multiple bandwidths

Uniform weights: $w_{\lambda} := 1 / |\Lambda|$

Number of permutations / wild bootstraps: $B_1 = B_2 = 2000$

JAX: runs on either CPU or GPU (significant speed improvements)

■ JAX GPU runs 100 times faster than Numpy CPU

mmdagg package: github.com/antoninschrab/mmdagg

```
from mmdagg import mmdagg          # X shape (m, d)
output = mmdagg(X, Y) # 0 or 1      # Y shape (n, d)
```

MMDAgg parameter-free user-friendly implementation

Radial basis function (RBF) kernel: $k_{\lambda}(\mathbf{x}, \mathbf{y}) := K\left(\left\|\frac{\mathbf{x} - \mathbf{y}}{\lambda}\right\|\right)$

Collection of bandwidths Λ : discretisation of the interval $[\lambda_{\min}, \lambda_{\max}]$ where λ_{\min} and λ_{\max} are the (robust) minimum and maximum of $\{\|\mathbf{x} - \mathbf{y}\| : \mathbf{x} \in \mathbb{X}_m, \mathbf{y} \in \mathbb{Y}_n\}$

Possible to aggregate several kernels each with multiple bandwidths

Uniform weights: $w_{\lambda} := 1 / |\Lambda|$

Number of permutations / wild bootstraps: $B_1 = B_2 = 2000$

JAX: runs on either CPU or GPU (significant speed improvements)

■ JAX GPU runs 100 times faster than Numpy CPU

mmdagg package: github.com/antoninschrab/mmdagg

```
from mmdagg import mmdagg          # X shape (m, d)
output = mmdagg(X, Y) # 0 or 1      # Y shape (n, d)
```

MMDAgg parameter-free user-friendly implementation

Radial basis function (RBF) kernel: $k_{\lambda}(\mathbf{x}, \mathbf{y}) := K\left(\left\|\frac{\mathbf{x} - \mathbf{y}}{\lambda}\right\|\right)$

Collection of bandwidths Λ : discretisation of the interval $[\lambda_{\min}, \lambda_{\max}]$ where λ_{\min} and λ_{\max} are the (robust) minimum and maximum of $\{\|\mathbf{x} - \mathbf{y}\| : \mathbf{x} \in \mathbb{X}_m, \mathbf{y} \in \mathbb{Y}_n\}$

Possible to aggregate several kernels each with multiple bandwidths

Uniform weights: $w_{\lambda} := 1 / |\Lambda|$

Number of permutations / wild bootstraps: $B_1 = B_2 = 2000$

JAX: runs on either CPU or GPU (significant speed improvements)

■ JAX GPU runs 100 times faster than Numpy CPU

mmdagg package: github.com/antoninschrab/mmdagg

```
from mmdagg import mmdagg          # X shape (m, d)
output = mmdagg(X, Y) # 0 or 1      # Y shape (n, d)
```

MMDAgg parameter-free user-friendly implementation

Radial basis function (RBF) kernel: $k_{\lambda}(\mathbf{x}, \mathbf{y}) := K\left(\left\|\frac{\mathbf{x} - \mathbf{y}}{\lambda}\right\|\right)$

Collection of bandwidths Λ : discretisation of the interval $[\lambda_{\min}, \lambda_{\max}]$ where λ_{\min} and λ_{\max} are the (robust) minimum and maximum of $\{\|\mathbf{x} - \mathbf{y}\| : \mathbf{x} \in \mathbb{X}_m, \mathbf{y} \in \mathbb{Y}_n\}$

Possible to aggregate several kernels each with multiple bandwidths

Uniform weights: $w_{\lambda} := 1 / |\Lambda|$

Number of permutations / wild bootstraps: $B_1 = B_2 = 2000$

JAX: runs on either CPU or GPU (significant speed improvements)

■ JAX GPU runs 100 times faster than Numpy CPU

mmdagg package: github.com/antoninschrab/mmdagg

```
from mmdagg import mmdagg                                # X shape (m, d)
output = mmdagg(X, Y) # 0 or 1                             # Y shape (n, d)
```

MMDAgg parameter-free user-friendly implementation

Radial basis function (RBF) kernel: $k_{\lambda}(\mathbf{x}, \mathbf{y}) := K\left(\left\|\frac{\mathbf{x} - \mathbf{y}}{\lambda}\right\|\right)$

Collection of bandwidths Λ : discretisation of the interval $[\lambda_{\min}, \lambda_{\max}]$ where λ_{\min} and λ_{\max} are the (robust) minimum and maximum of $\{\|\mathbf{x} - \mathbf{y}\| : \mathbf{x} \in \mathbb{X}_m, \mathbf{y} \in \mathbb{Y}_n\}$

Possible to aggregate several kernels each with multiple bandwidths

Uniform weights: $w_{\lambda} := 1 / |\Lambda|$

Number of permutations / wild bootstraps: $B_1 = B_2 = 2000$

JAX: runs on either CPU or GPU (significant speed improvements)

■ JAX GPU runs 100 times faster than Numpy CPU

mmdagg package: github.com/antoninschrab/mmdagg

```
from mmdagg import mmdagg                                # X shape (m, d)
output = mmdagg(X, Y) # 0 or 1                             # Y shape (n, d)
```

MMDAgg parameter-free user-friendly implementation

Radial basis function (RBF) kernel: $k_{\lambda}(\mathbf{x}, \mathbf{y}) := K\left(\left\|\frac{\mathbf{x} - \mathbf{y}}{\lambda}\right\|\right)$

Collection of bandwidths Λ : discretisation of the interval $[\lambda_{\min}, \lambda_{\max}]$ where λ_{\min} and λ_{\max} are the (robust) minimum and maximum of $\{\|\mathbf{x} - \mathbf{y}\| : \mathbf{x} \in \mathbb{X}_m, \mathbf{y} \in \mathbb{Y}_n\}$

Possible to aggregate several kernels each with multiple bandwidths

Uniform weights: $w_{\lambda} := 1 / |\Lambda|$

Number of permutations / wild bootstraps: $B_1 = B_2 = 2000$

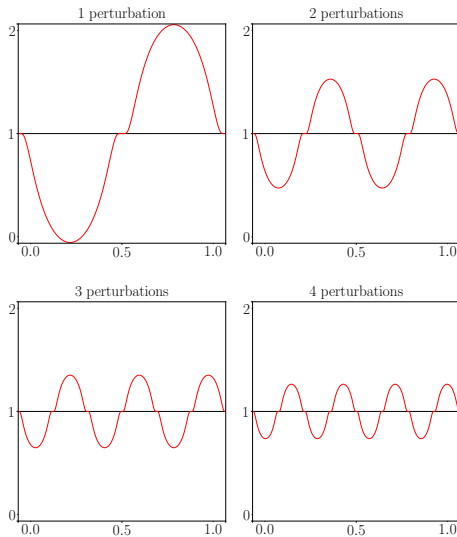
JAX: runs on either CPU or GPU (significant speed improvements)

■ JAX GPU runs 100 times faster than Numpy CPU

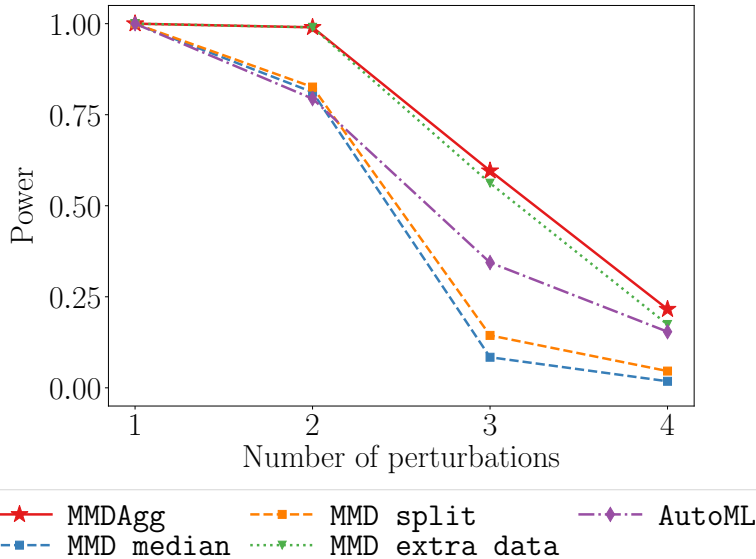
mmdagg package: github.com/antoninschrab/mmdagg

```
from mmdagg import mmdagg                # X shape (m, d)
output = mmdagg(X, Y) # 0 or 1            # Y shape (n, d)
```

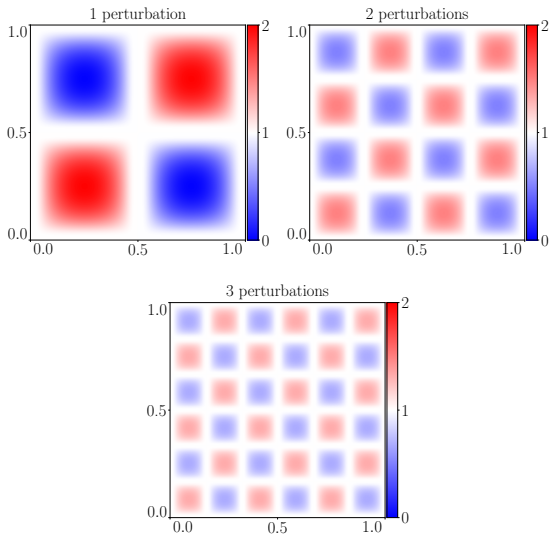
Experiment on perturbed uniform $d = 1$



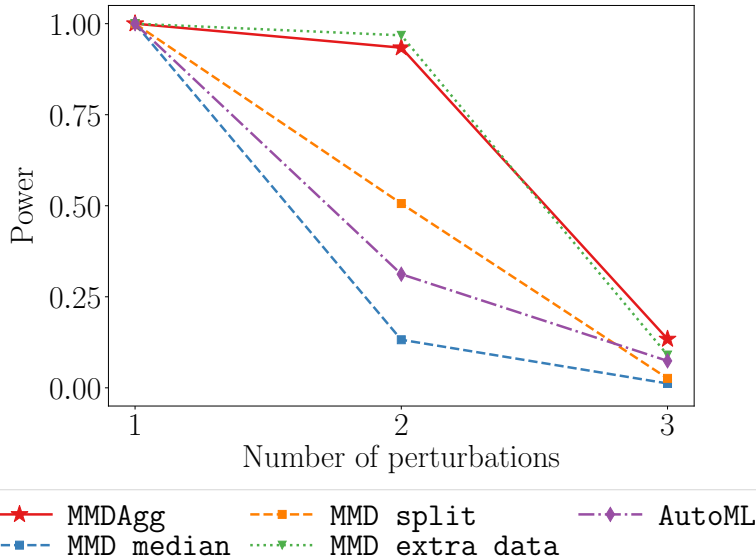
Experiment on perturbed uniform $d = 1$



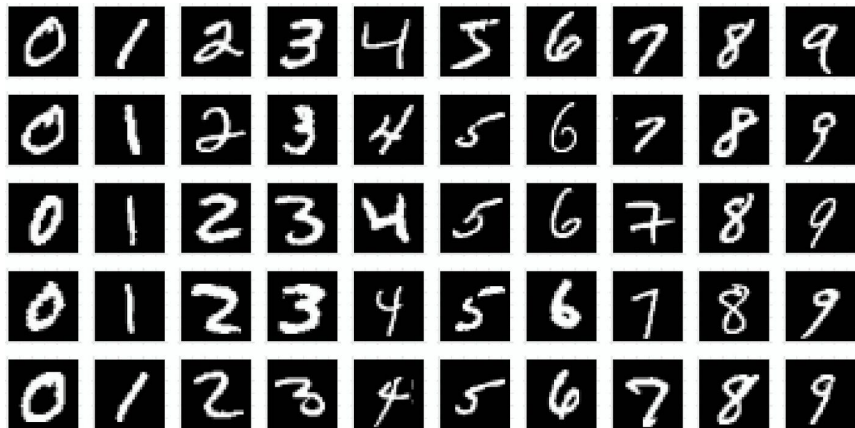
Experiment on perturbed uniform $d = 2$



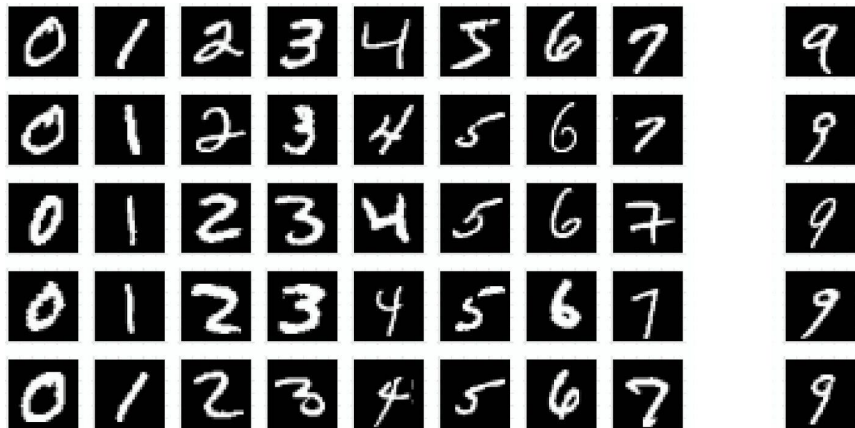
Experiment on perturbed uniform $d = 2$



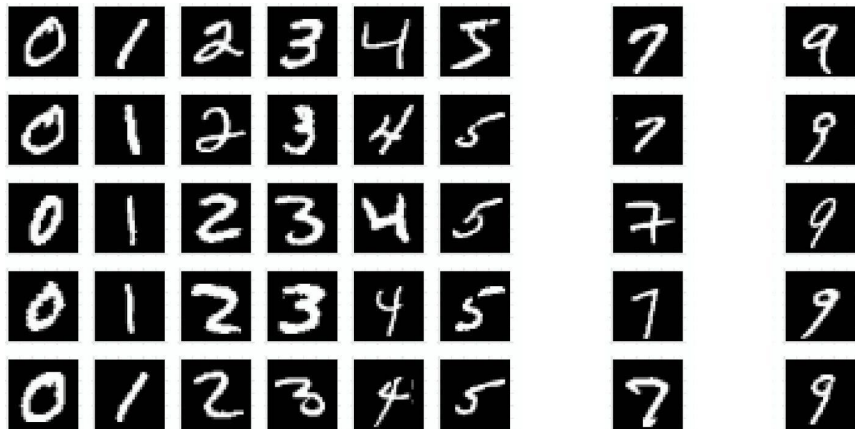
Experiment on MNIST digits



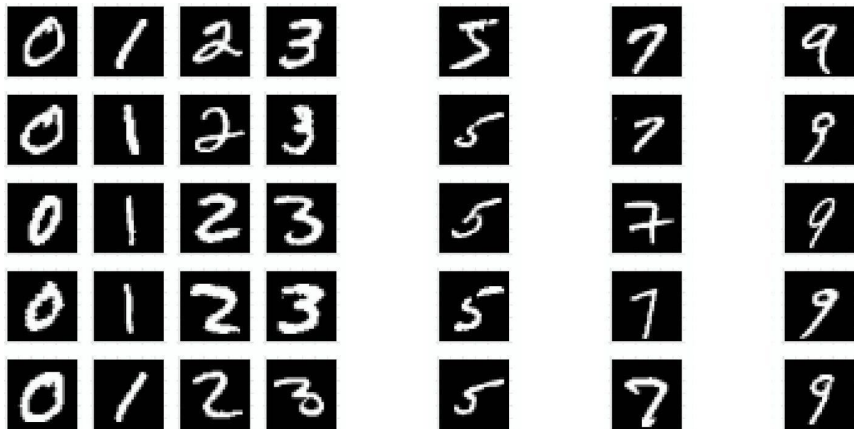
Experiment on MNIST digits



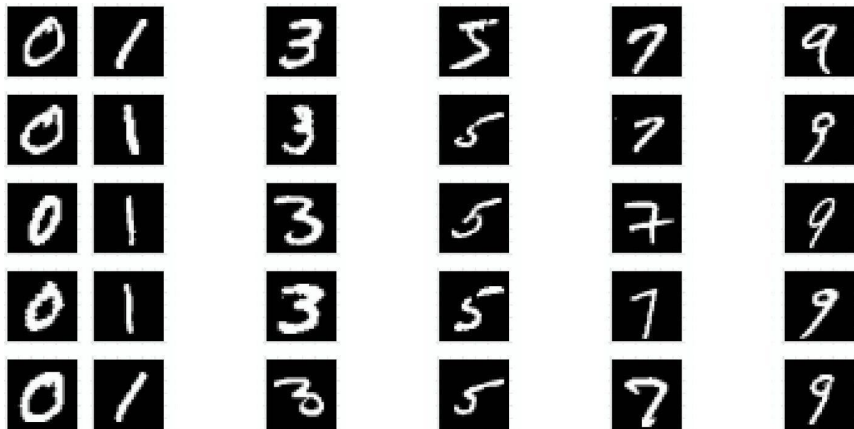
Experiment on MNIST digits



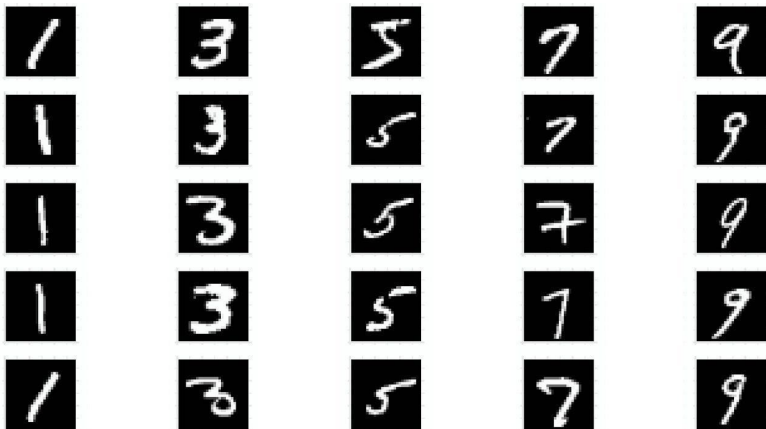
Experiment on MNIST digits



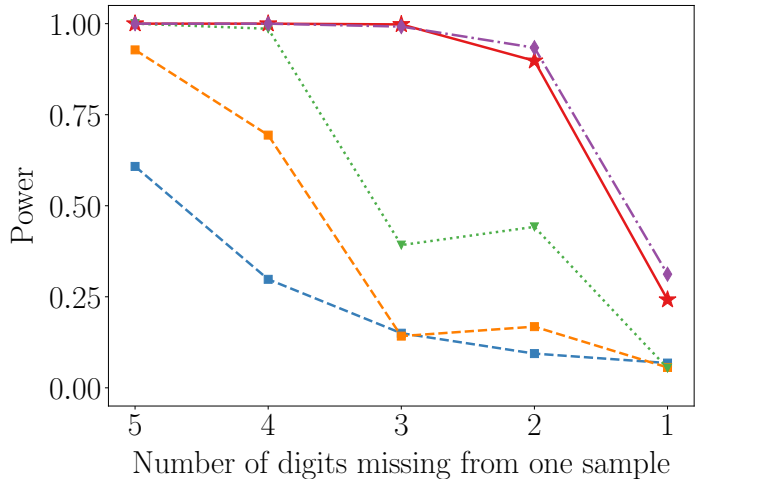
Experiment on MNIST digits



Experiment on MNIST digits

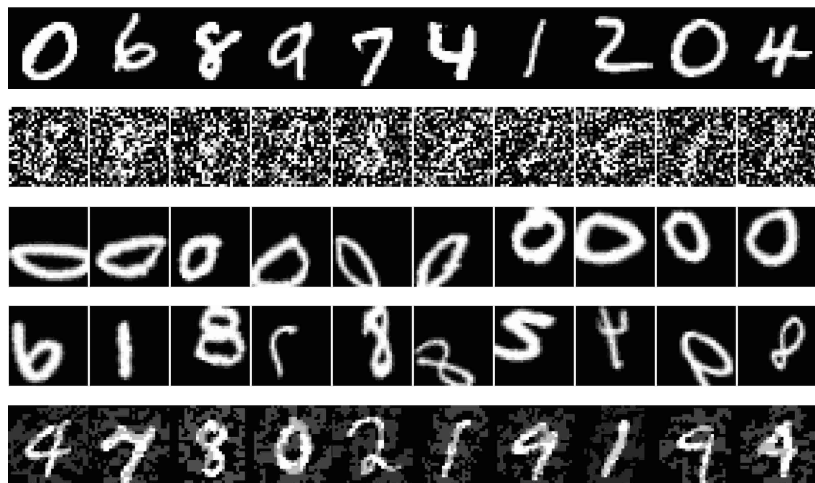


Experiment on MNIST digits



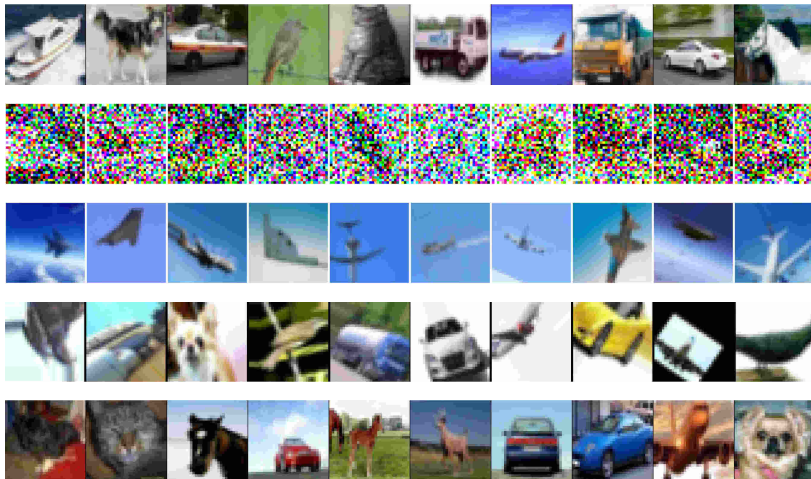
Experiment on image shifts on MNIST & CIFAR-10

Failing Loudly Benchmark: Rabanser et al., 2019

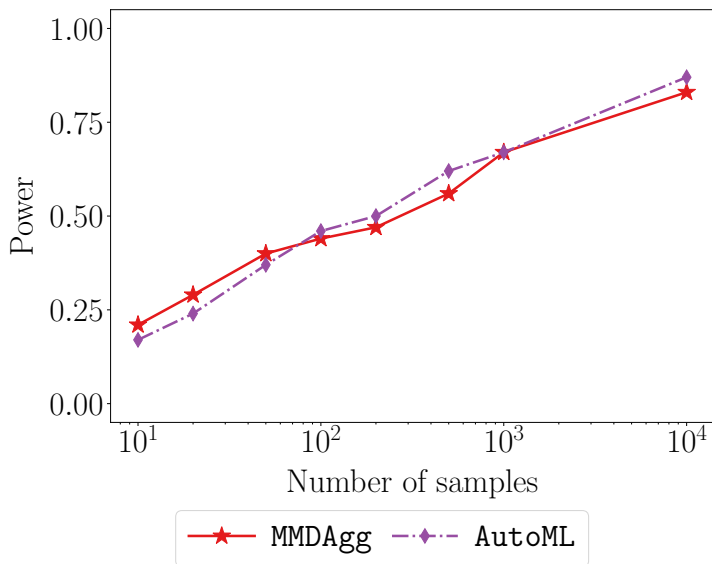


Experiment on image shifts on MNIST & CIFAR-10

Failing Loudly Benchmark: Rabanser et al., 2019



Experiment on image shifts on MNIST & CIFAR-10



MMD kernel choice without data splitting

MMD Aggregated Two-Sample Test (JMLR 2023):



Code:

<https://github.com/antoninschrab/mmdagg-paper>

Research support

Work supported by:

The Gatsby Charitable Foundation



Google Deepmind



Questions?

