Compress Then Test: Powerful Kernel Testing in Near-linear Time



Carles Domingo-Enrich







Massachusetts Institute of Fechnology



Raaz Dwivedi





Lester Mackey



JSM 2023

Two-sample testing

Test null hypothesis: \mathcal{H}_0 : $\mathbb{P} = \mathbb{Q}$

with $\mathbf{X}_n = (X_i)_{1 \le i \le n}$ i.i.d. sample of \mathbb{P} $\mathbf{Y}_n = (Y_i)_{1 \le i \le n}$ i.i.d. sample of \mathbb{Q}



• An alternate representation of \mathcal{H}_0 : $\mathbb{P} = \mathbb{Q}$

 $\mathbb{E}_{X \sim \mathbb{P}} f(X) = \mathbb{E}_{Y \sim \mathbb{Q}} f(Y) \text{ for all } f \in \mathcal{F} \text{ (a class of test functions)}$



- An alternate representation of \mathcal{H}_0 : $\mathbb{P} = \mathbb{Q}$
- Parametric test: When \mathcal{F} is finite-dimensional (e.g., $\mathcal{F} = \{X, X^2\}$ for testing Gaussian \mathbb{P}, \mathbb{Q})

 $\mathbb{E}_{X \sim \mathbb{P}} f(X) = \mathbb{E}_{Y \sim \mathbb{Q}} f(Y)$ for all $f \in \mathcal{F}$ (a class of test functions)



• An alternate representation of $\mathcal{H}_0: \mathbb{P} = \mathbb{Q}$

- Parametric test: When \mathcal{F} is finite-dimensional (e.g., $\mathcal{F} = \{X, X^2\}$ for testing Gaussian \mathbb{P}, \mathbb{Q})
- Non-parametric test: When \mathcal{F} is infinite-dimensional (e.g., $\mathcal{F} =$ unit ball of a reproducing kernel Hilbert space)

 $\mathbb{E}_{X \sim \mathbb{P}} f(X) = \mathbb{E}_{Y \sim \mathbb{Q}} f(Y) \text{ for all } f \in \mathcal{F} \text{ (a class of test functions)}$



• A symmetric function $\mathbf{k} : \mathbb{R}^d \times \mathbb{R}^d \to \mathbb{R}$, i.e., $\mathbf{k}(x, y) = \mathbf{k}(y, x)$ which is positive semidefinite, is called a reproducing kernel

• A symmetric function $\mathbf{k} : \mathbb{R}^d \times \mathbb{R}^d \to \mathbb{R}$, i.e., $\mathbf{k}(x, y) = \mathbf{k}(y, x)$







which is positive semidefinite, is called a reproducing kernel

- A symmetric function $\mathbf{k} : \mathbb{R}^d \times \mathbb{R}^d \to \mathbb{R}$, i.e., $\mathbf{k}(x, y) = \mathbf{k}(y, x)$ which is positive semidefinite, is called a reproducing kernel
- Reproducing kernel Hilbert space = Closure of linear span of $\{\mathbf{k}(x, \cdot)\}$





ullet Reparametrized via maximum mean discrepancy (MMD) of a kernel ${f k}$

$$\mathscr{H}_{0}: \mathrm{MMD}_{\mathbf{k}}(\mathbb{P}, \mathbb{Q}) \stackrel{\Delta}{=} \sup_{\|f\|_{\mathbf{k}} \leq 1} \mathbb{E}_{X \sim \mathbb{P}} f(X)$$

 $X) - \mathbb{E}_{Y \sim \mathbb{Q}} f(Y) = 0$



- MMD expression:

$\mathscr{H}_{0}: \mathrm{MMD}_{\mathbf{k}}^{2}(\mathbb{P}, \mathbb{Q}) \stackrel{\Delta}{=} \mathbb{E}_{X, X' \sim \mathbb{P}} \mathbf{k}(X, X') + \mathbb{E}_{Y, Y' \sim \mathbb{P}} \mathbf{k}(Y, Y') - 2\mathbb{E}_{X \sim \mathbb{P}, Y \sim \mathbb{Q}} \mathbf{k}(X, Y) = 0$



- MMD expression:
- Test statistic:

 $MMD_{\mathbf{k}}^{2}(\mathbf{X}_{n},\mathbf{Y}_{n}) = \frac{\sum_{i,i'=1}^{n} \mathbf{k}(X_{i},X_{i'})}{n^{2}} + \frac{\sum_{j,j'=1}^{n} \mathbf{k}(Y_{j},Y_{j'})}{n^{2}} - \frac{2\sum_{i=1}^{n} \sum_{j=1}^{n} \mathbf{k}(X_{i},Y_{j})}{n^{2}}$

$\mathscr{H}_{0}: \mathrm{MMD}_{\mathbf{k}}^{2}(\mathbb{P}, \mathbb{Q}) \stackrel{\Delta}{=} \mathbb{E}_{X, X' \sim \mathbb{P}} \mathbf{k}(X, X') + \mathbb{E}_{Y, Y' \sim \mathbb{P}} \mathbf{k}(Y, Y') - 2\mathbb{E}_{X \sim \mathbb{P}, Y \sim \mathbb{O}} \mathbf{k}(X, Y) = 0$



- MMD expression:
- Test statistic:

$$\mathbf{MMD}_{\mathbf{k}}^{2}(\mathbf{X}_{n}, \mathbf{Y}_{n}) = \frac{\sum_{i,i'=1}^{n} \mathbf{k}(X_{i}, X_{i'})}{n^{2}}$$

• **Test:** If $MMD_k^2(\mathbf{X}_n, \mathbf{Y}_n) < t_\alpha$ accept the null; else reject

$\mathscr{H}_{0}: \mathrm{MMD}_{\mathbf{k}}^{2}(\mathbb{P}, \mathbb{Q}) \stackrel{\Delta}{=} \mathbb{E}_{X, X' \sim \mathbb{P}} \mathbf{k}(X, X') + \mathbb{E}_{Y, Y' \sim \mathbb{P}} \mathbf{k}(Y, Y') - 2\mathbb{E}_{X \sim \mathbb{P}, Y \sim \mathbb{Q}} \mathbf{k}(X, Y) = 0$

 $+ \frac{\sum_{j,j'=1}^{n} \mathbf{k}(Y_{j}, Y_{j'})}{n^{2}} - \frac{2 \sum_{i=1}^{n} \sum_{j=1}^{n} \mathbf{k}(X_{i}, Y_{j})}{n^{2}}$



- MMD expression:
- Test statistic:

$$\mathrm{MMD}_{\mathbf{k}}^{2}(\mathbf{X}_{n},\mathbf{Y}_{n}) = \frac{\sum_{i,i'=1}^{n} \mathbf{k}(X_{i},X_{i'})}{n^{2}} + \frac{1}{n^{2}}$$

• **Test:** If $MMD_k^2(\mathbf{X}_n, \mathbf{Y}_n) < t_{\alpha}$ accept the null; else reject

$\mathscr{H}_{0}: \mathrm{MMD}_{\mathbf{k}}^{2}(\mathbb{P}, \mathbb{Q}) \stackrel{\Delta}{=} \mathbb{E}_{X, X' \sim \mathbb{P}} \mathbf{k}(X, X') + \mathbb{E}_{Y, Y' \sim \mathbb{P}} \mathbf{k}(Y, Y') - 2\mathbb{E}_{X \sim \mathbb{P}, Y \sim \mathbb{Q}} \mathbf{k}(X, Y) = 0$





Type 1 error and validity of test

• Valid test: Threshold t_{α} chosen such that

 $\mathcal{H}_0: \mathbb{P} = \mathbb{Q}$

Type 1 error = $\operatorname{Prob}_{\mathcal{H}_0} \left| \operatorname{MMD}_{\mathbf{k}}^2(\mathbf{X}_n, \mathbf{Y}_n) \ge t_{\alpha} \right| \le \alpha$





of test statistics under the null

• Asymptotic tests: Use quantiles from the asymptotic distribution



Asymptotic tests: Use quantile of test statistics under the null

 n² runtime

• Asymptotic tests: Use quantiles from the asymptotic distribution



Asymptotic tests: Use quantile of test statistics under the null

 n² runtime

inexact with finite samples; asymptotics might not be Gaussian

• Asymptotic tests: Use quantiles from the asymptotic distribution



- of test statistics under the null $\sqrt{n^2}$ runtime
- of test statistics under the null

• Asymptotic tests: Use quantiles from the asymptotic distribution

inexact with finite samples; asymptotics might not be Gaussian

• **Permutation tests:** Use permutations to simulate the distribution



- of test statistics under the null $\sqrt{n^2}$ runtime
- of test statistics under the null vert type 1 error control in finite samples

• Asymptotic tests: Use quantiles from the asymptotic distribution

inexact with finite samples; asymptotics might not be Gaussian

• Permutation tests: Use permutations to simulate the distribution



- of test statistics under the null $\sqrt{n^2}$ runtime
- of test statistics under the null exact type 1 error control in finite samples

• Asymptotic tests: Use quantiles from the asymptotic distribution

inexact with finite samples; asymptotics might not be Gaussian

• Permutation tests: Use permutations to simulate the distribution

 \bullet $n^2 \mathscr{L}$ runtime, when doing \mathscr{L} permutations (typically $\mathscr{L} \sim 100$ s)



Power of the test and MMD separation



Power of the test and MMD separation

• Power of the test = 1 – Type 2 error $\mathscr{H}_1 : \mathbb{P} \neq \mathbb{Q}$ $= \operatorname{Prob}_{\mathcal{H}_{1}} \left[\operatorname{MMD}_{\mathbf{k}}^{2}(\mathbf{X}_{n}, \mathbf{Y}_{n}) \geq t_{\alpha} \right]$





Power of the test and MMD separation

• Power of the test = 1 - Type 2 error

• MMD separation rate_n = Minimum value of $MMD_{k}(\mathbb{P}, \mathbb{Q})$ such that power of the test is bounded away from 0

 $\mathcal{H}_1: \mathbb{P} \neq \mathbb{Q}$ $= \operatorname{Prob}_{\mathcal{H}_{1}} \left[\operatorname{MMD}_{\mathbf{k}}^{2}(\mathbf{X}_{n}, \mathbf{Y}_{n}) \geq t_{\alpha} \right]$





MMD separation Runtime Test name

Complete MMD

(Gretton et al., 2012a)

 $n^{-\frac{1}{2}}$

 n^2



Our goal: Speed up the test runtime while maintaining

- 1. validity for all sample sizes (exact type 1 error)
- error)

2. small MMD separation rate (high power, low type 2

MMD separation Test name Runtime



B = Number of blocks used in Block MMD test

Test name

MMD separation Runtime

Complete MMD (Gretton et al., 2012a)	$n^{-\frac{1}{2}}$	n^2
Block MMD (Zaremba et al., 2013)	$(Bn)^{-\frac{1}{4}}$	Bn
Incomplete MMD (Yamada et al., 2019)	$e^{-\frac{1}{4}}$	ŀ
	B = Num $\ell = Num$	ber of b
	test	

locks used in Block MMD test rdered index pairs in Incomplete MMD





Prior strategies to speed up testing





Prior strategies to speed up testing: Faster runtimes, but only at the loss of power



Gaussian (mean separation = 0.012, $n = 4^9$)

Prior strategies to speed up testing: Faster runtimes, but only at the loss of power



Gaussian (mean separation = 0.012, $n = 4^9$)



Our solution: Compress Then Test Best of three worlds: Fast + Always valid + High power



2 5 20 10

Gaussian (mean separation = 0.012, $n = 4^9$)

50 100 200 500 1k 2k 5k Total computation time (s)



Key idea: Compress the data and then test

Compress \mathbf{X}_n to $\hat{\mathbf{X}}_m$ of size m

Random
subsamplingKT-Compress
(Shetty-Dwivedi-Mackey,
2022)MMD_k(X_n, \hat{X}_m) $O(n^{-1/4})$ $\tilde{O}(n^{-1/2})$

$$n = \sqrt{n}$$
 with small $MMD_k^2(\mathbf{X}_n, \hat{\mathbf{X}}_m)$





Key idea: Compress non-uniformly for high quality Compress X_n to \hat{X}_m of size $m = \sqrt{n}$ with small MMD_k²(X_n, \hat{X}_m) Random sub-sampling 15.0 **KT-Compress** 88





KT-Compress = Kernel thinning (KT) + Compress



KT-Compress = Kernel thinning (KT) + Compress KT: Take n^2 time to output *m* points with $MMD_{k}(\mathbf{X}_{n}, \hat{\mathbf{X}}_{m}) = \tilde{O}(1/m)$



KT-Compress = Kernel thinning (KT) + Compress KT: Take n^2 time to output *m* points with $MMD_{k}(\mathbf{X}_{n}, \hat{\mathbf{X}}_{m}) = \tilde{O}(1/m)$

and an <u>arbitrary</u> kernel (Gaussian, Matérn, ...)

- Works with an <u>arbitrary</u> sequence of points (IID, MCMC, ...)



KT-Compress = Kernel thinning (KT) + Compress KT: Take n^2 time to output *m* points with $MMD_{k}(\mathbf{X}_{n}, \hat{\mathbf{X}}_{m}) = \tilde{O}(1/m)$

and an <u>arbitrary</u> kernel (Gaussian, Matérn, ...)

Kernel thinning, [Dwivedi & Mackey '21]



- Works with an <u>arbitrary</u> sequence of points (IID, MCMC, ...)





KT-Compress = Kernel thinning (KT) + Compress



KT-Compress = Kernel thinning (KT) + Compress Compress improves KT's runtime from $\Theta(n^2)$ to $\tilde{O}(n)$ while making MMD error worse by at most a factor of 4!

Also, Compress provides a quadratic speedup to the runtime of <u>generic</u> thinning algorithms!



KT-Compress = Kernel thinning (KT) + Compress Compress improves KT's runtime from $\Theta(n^2)$ to $\tilde{O}(n)$ while making MMD error worse by at most a factor of 4!

Also, Compress provides a quadratic speedup to the runtime of <u>generic</u> thinning algorithms!



Distribution compression in nearlinear time, [Shetty, Dwivedi & Mackey '22]



Kernel thinning: Algorithm overview

Recursive kernel halving: $n \to \frac{n}{2} \to \frac{n}{4} \to \dots \to \sqrt{n}$

+ Post-processing step



Kernel thinning: Algorithm overview

Recursive kernel halving: $n \to \frac{n}{2} \to \frac{n}{4} \to \dots \to \sqrt{n}$



+ Post-processing step S_{in} $S^{1,2}$ \$2,3 $S^{2,4}$ $S^{3,6}$ $S^{3,8}$ $S^{3,7}$ $S^{m,4}$ $S^{m,5}$ $S^{m,6}$ $S^{m,7}$ $S^{m,8}$ \cdots $S^{m,2^m}$



Kernel halving = Sequential non-uniform selection





 $s^{(m,1)}$ (n points)



Kernel halving = Sequential non-uniform selection



Sequential selection of^ma point out of paired points, using biased coin flips toward the point with smaller MMD error to input points $S^{(m,1)}$ (n points)





Compress: Recursive strategy to reduce runtime

Thinning algorithm with n^2 runtime and \rightarrow error e(n)



Algorithm 1: COMPRESS

if $n = 4^{g}$ then return S_{in} else

Partition S_{in} into four arbitrary subsequences $\{S_i\}_{i=1}^4$ each of size n/4for i = 1, 2, 3, 4 do

 $\widetilde{S}_i \leftarrow \text{COMPRESS}(S_i, \text{HALVE}, \mathfrak{g})$ // return coresets of size $2^{\mathfrak{g}} \cdot \sqrt{\frac{n}{4}}$ end

 $\widetilde{S} \leftarrow \text{CONCATENATE}(\widetilde{S}_1, \widetilde{S}_2, \widetilde{S}_3, \widetilde{S}_4)$ return HALVE(S)

end

 \sqrt{n} -thinning algorithm with \rightarrow $n \log^3 n$ runtime and error 4e(n)

Input: halving algorithm HALVE, oversampling parameter g, point sequence S_{in} of size n

// coreset of size $2 \cdot 2^{\mathfrak{g}} \cdot \sqrt{n}$ // coreset of size $2^{\varrho}\sqrt{n}$





1. **KT-Compress:** $\mathbf{X}_n \to \hat{\mathbf{X}}_m \& \mathbf{Y}_n \to \hat{\mathbf{Y}}_m$ with size $m = 2^{\mathfrak{g}} \sqrt{n}$





2. Test statistic: $MMD_k^2(\mathbf{X}_n, \mathbf{Y}_n) \rightarrow MMD_k^2(\hat{\mathbf{X}}_m, \hat{\mathbf{Y}}_m)$



- 1. **KT-Compress:** $X_n \to \hat{X}_m \& Y_n \to \hat{Y}_m$ with size $m = 2^{\mathfrak{g}} \sqrt{n}$
- 2. Test statistic: $MMD_k^2(X_n, Y_n) \rightarrow MMD_k^2(\hat{X}_m, \hat{Y}_m)$



3. Threshold: Via cheap permutations, group data into $s \ll \sqrt{n}$ bins, sample \mathscr{L} permutations of [s], and permute the s bins



1. **KT-Compress:** $\mathbf{X}_n \to \hat{\mathbf{X}}_m \& \mathbf{Y}_n \to \hat{\mathbf{Y}}_m$ with size $m = 2^{\mathfrak{g}} \sqrt{n}$

2. Test statistic: $MMD_k^2(\mathbf{X}_n, \mathbf{Y}_n) \rightarrow MMD_k^2(\hat{\mathbf{X}}_m, \hat{\mathbf{Y}}_m)$

CTT runtime: $\tilde{O}(n) + O(s^2 \mathscr{L})$ if $\mathfrak{q} = \log \log n$

Complete MMD test runtime: $O(n^2) + O(n^2\mathscr{L})$



3. Threshold: Via cheap permutations, group data into $s \ll \sqrt{n}$ bins, sample \mathscr{L} permutations of [s], and permute the s bins



Test name	MMD separation	Runtin	
CTT (ours)	$\frac{R_{\mathbf{k},\mathbf{k}'}(\mathbb{P},n,\delta,\mathfrak{g})}{2\mathfrak{g}\sqrt{n}}+n^{-\frac{1}{2}}$	4 ^g n 10g	
Complete MMD (Gretton et al., 2012a)	$n^{-\frac{1}{2}}$	n^2	
Block MMD (Zaremba et al., 2013)	$(Bn)^{-\frac{1}{4}}$	Bn	
Incomplete MMD (Yamada et al., 2019)	$\mathcal{C}^{-\frac{1}{4}}$	l	

B = Number of blocks used in Block MMD test $\delta =$ Failure probability for CTT guarantees $\mathbf{k}' =$ auxiliary kernel used by KT-Compress

ne

g n



- ℓ = Number of ordered index pairs in Incomplete MMD test



Test name	MMD separation	Runtime	Tails of \mathbb{P} on \mathbb{R}^d	Choice of k'	$R_{\mathbf{k},\mathbf{k}'}(\mathbb{P},n)$
CTT (ours)	$\frac{R_{\mathbf{k},\mathbf{k}'}(\mathbb{P},n,\delta,\mathfrak{g})}{2\mathfrak{g}\sqrt{n}}+n^{-\frac{1}{2}}$	$4^{\mathfrak{g}}n\log n$	Compact	Compact k _{rt}	$\left(\log\frac{n}{\delta}\right)$
Complete MMD (Gretton et al., 2012a)	$n^{-\frac{1}{2}}$	n^2	Subexp	Analytic k	$\left(\log\frac{n}{\delta}\right)$
Block MMD (Zaremba et al., 2013)	$(Bn)^{-\frac{1}{4}}$	Bn	Subexp	Subexp k _{rt}	$c_{m,\delta}(\log \frac{n}{\delta})$
Incomplete MMD (Yamada et al., 2019)	$\mathcal{C}^{-\frac{1}{4}}$	ŀ	ho-Heavy-tail	$ ho$ -Heavy-tail ${f k}_{ m rt}$	$\left(\frac{n}{\delta}\right)^{\frac{d}{2\rho}} \left(\log\frac{n}{\delta}\right)$
$B = \mathcal{E} = \delta$ $\delta = \mathbf{k'} = $	 Number of block Number of order Failure probabilit auxiliary kernel u 	ks used in B red index p sy for CTT g used by KT-(olock MMD test airs in Incomple uarantees Compress & k (<i>x</i>	te MMD test , y)= $\int \mathbf{k}_{rt}(x, z) \mathbf{l}$	$\mathbf{k}_{\mathrm{rt}}(z,y)dz$





Time-power tradeoff on Gaussian data



Gaussian (mean separation = 0.012, $n = 4^9$)



inear Time [Domingo-Enrich, Dwivedi & Mackey '23]

- 1. Quadratic \rightarrow Near-linear time tests using KT-Compress
- 2. Provide up to 200x speed-up (1 hour \rightarrow 20 sec)
- 3. Maintain level <u>and</u> power provably
- 4. Work with kernel approximations (Low-rank CTT)
- 5. Applicable for kernel selection (Aggregated CTT)

Compress Then Test: Powerful Kernel Testing in Near-





Kernel Thinning (KT)

KT-Compress

Compress then Test (CTT)



n points $\rightarrow \sqrt{n}$ points with $\tilde{O}(n^{-1/2})$ error in n^2 time

 $\tilde{O}(n)$ runtime version of KT with error worse by a factor of 4

\sim Powerful O(n) runtime kernel tests using KT-Compress









Appendix

Low-rank CTT vs random Fourier features based test





Aggregated CTT for kernel selection



Total computation time (s)





Abhishek Shetty

Carles

Domingo

Enrich

Lester Mackey

- 1. Kernel Thinning, COLT 2021
- 2. Generalized Kernel Thinning, ICLR 2022

 $\rightarrow \sqrt{n}$ points with $\Theta(n^{-1/4})$ error (Monte Carlo rate)

Powerful always valid two-sample test in $\tilde{O}(n)$ time

pip install goodpoints

3. Distribution Compression In Near-linear Time, ICLR 2022

4. Compress Then Test: Powerful Kernel Testing in Near-Linear Time, AISTATS 2023









Compress++: Recursive strategy to reduce runtime

Thinning algorithm with n^2 runtime and error e(n)



Input: halving algorithm HALVE, oversampling parameter g, point sequence S_{in} of size n if $n = 4^{g}$ then return S_{in} else

Compress++

Partition S_{in} into four arbitrary subsequences $\{S_i\}_{i=1}^4$ each of size n/4for i = 1, 2, 3, 4 do

end

 $\widetilde{S} \leftarrow \text{CONCATENATE}(\widetilde{S}_1, \widetilde{S}_2, \widetilde{S}_3, \widetilde{S}_4)$ return HALVE(S)

end

 \sqrt{n} -thinning algorithm with $n \log^3 n$ runtime and error 4e(n)

 $\widetilde{S}_i \leftarrow \text{COMPRESS}(S_i, \text{HALVE}, \mathfrak{g})$ // return coresets of size $2^{\mathfrak{g}} \cdot \sqrt{\frac{n}{4}}$

// coreset of size $2 \cdot 2^{\mathfrak{g}} \cdot \sqrt{n}$ // coreset of size $2^{\varrho}\sqrt{n}$