Comparing distributions: ℓ_1 geometry improves kernel two-sample testing

Overview

Problem: Are two sets of observations drawn from the same distribution?

Contributions:

- We exhibit a family of *L^p*-based metrics which metrize the weak convergence.
- We derive linear-time, nonparametric, a.s consistent *L*¹ based two sample tests.
- We show L^1 geometry provides better power than its L^2 counterpart.
- We maximize a lower bound on the test power and learn distinguishing features between distributions.



Theorem: Let **k** a characteric and bounded kernel. For all $p \geq 1$,

$$d_{L^p,\mu}(\boldsymbol{P},\boldsymbol{Q}) \coloneqq \left(\int_t |\boldsymbol{\mu}_{\boldsymbol{P}}(t) - \boldsymbol{\mu}_{\boldsymbol{Q}}(t)|^p d\boldsymbol{\Gamma}(t)\right)^{1/p}$$

where $\mu_P(t) := \int_t k(x, t) dP(x)$ is a metric which metrize the weak convergence.

Sketch of proof:

- Integral operator: $T_{k_{\sigma}}: f \in L_2^{d\Gamma}(\mathbb{R}^d) \to \int_{\mathbb{R}^d} k(x, .) f(x) dI$
- Unit Ball of $L^{d\Gamma}_{\infty}(\mathbb{R}^d)$: $B^{d\Gamma}_{\infty} \coloneqq \{f: \sup |f(x)| \le 1 \text{ a.s}\}$
- **IPM** formulation: $d_{L^p,\mu}(P,Q) = \sup \{E_P(f(X)) E_Q(f(Y))\}$ $f \in T_k(B^{d\Gamma}_{\infty})$

Mean Embedding test

• Test: H_0 : $\mathbf{P} = \mathbf{Q}$ vs H_1 : $\mathbf{P} \neq \mathbf{Q}$: • Samples: $X \coloneqq \{x_i\}_{i=1}^n \sim P$ and $Y \coloneqq \{y_i\}_{i=1}^n \sim Q$ • Empirical ME: $\mu_X(T) \coloneqq \frac{1}{n} \sum_{i=1}^n k_\sigma(x_i, T)$ • k_{σ} the Gaussian kernel of width σ • Test locations: $\{T_i\}_{i=1}^{J} \sim \Gamma$ • Test statistic:

$$\left(\widehat{d}_{\ell_p,\mu}(X,Y)\right)^p \coloneqq n^{\frac{p}{2}} \sum_{i=1}^J |\mu_X(T_i) - \mu_Y(T_i)|^p$$



Test of level α : Compute $\left(\widehat{d}_{\ell_p,\mu}(X,Y)\right)^{\prime}$ and reject H_0 if $\left(\widehat{d}_{\ell_{p},\mu}(X,Y)\right)^{p} > T_{\alpha,p} = 1 - \alpha$ quantile of the null distribution.

Why $\ell_1 \gg \ell_2$?

Definition: (Analytic kernel). *A positive definite kernel* k is analytic if for all $x \in \mathbb{R}^d$, the feature map k(x, .) is an analytic function on \mathbb{R}^d .

Proposition: Let $\delta > 0$. Under the null hypothesis H_0 , almost surely there exists $N \ge 1$, such that for all $n \geq N$, with a probability of 1- δ :

$$\left(\widehat{d}_{\ell_{2},\mu}(X,Y)\right)^{2} > T_{\alpha,2} \implies \widehat{d}_{\ell_{1},\mu}(X,Y) > T_{\alpha,\mu}$$







Regularization: To obtain a lower bound, we consider the regularized statistic

Normalized Tests

Remark: Under H_0 , $\hat{d}_{\ell_1,\mu}(X, Y)$ converge to a sum of correlated Nakagami variables.

Normalized Mean Embedding (ME) Test:

 $L1-ME[X,Y] \coloneqq ||\sqrt{n}\Sigma_n^{-2}S_n||_1$ • $S_n \coloneqq \frac{1}{n} \sum_{i=1}^n Z_X^i - Z_Y^i$ • $\Sigma_n \coloneqq \widehat{cov}(Z_X) + \widehat{cov}(Z_Y)$ • $Z_X^i \coloneqq (k_\sigma(x_i, T_1), \dots, k_\sigma(x_i, T_J))$

Proposition: Under *H*₀, **L1-ME**[*X*, *Y*] is a.s asymptotically distributed as a sum of J i.i.d Nakagami variables of parameter $m = \frac{1}{2}$ and $\varpi = \frac{1}{2}$.

Normalized Smooth Characteristic Function (SCF) Test:

 $L1-SCF[X,Y] \coloneqq ||\sqrt{n}\Sigma_n^{-\overline{2}}S_n||_1$

• $Z_X^i \coloneqq (cos(x_i^T T_1)f(x_i), sin(x_i^T T_1)f(x_i), \dots, sin(x_i^T T_J)f(x_i))$ f is the inverse Fourier transform of k_{σ} .

Optimization Procedure

L1-ME[X, Y] := $||\sqrt{n} (\Sigma_n + \gamma_n)^{-1/2} S_n||_1$ • $\gamma_n \rightarrow 0$

Proposition: The test power $P(L1 - ME[X, Y] > \epsilon)$ of the the L1-ME test satisifies $P(L1-ME[X,Y] > \epsilon) \ge L(\lambda_n)$ where $L(\lambda_n)$ is an increasing function of λ_n and goes to 1 when *n* goes to infinity.

• $\lambda_n \coloneqq ||\sqrt{n} \Sigma^2 S||_1$ is the population counterpart of L1-ME[X, Y].

Optimization Procedure:

• Optimize $\{T_i\}_{i=1}^J, \sigma = \operatorname{argmax} L(\lambda_n)$

• Estimation of λ_n on a separate training set.





L1	-opt
I 1	ani

power.

Meyer Scetbon, Gaël Varoquaux Inria, Université Paris-Saclay

Informative Features

Contour plot of L1–ME[X, Y] as a function of T_2 with J = 2 and T_1 fixed.

- $P \sim N([0, 0], I_2)$
- $Q \sim N([0, 1], I_2)$
- L1–ME[*X*, *Y*] detects the differences.

Test Power: Synthetic Problems

• Test Power vs. Sample Size



t-ME, L1-opt-SCF: Proposed Methods L1-grid-ME, L1-grid-SCF: Random settings • ME-full, SCF-full: Optimized ℓ_2 -based methods



L1-opt-SCF

 L1-grid-SCF ME-full

SCF-full

--- MMD-lin MMD-quad

• MMD-quad, MMD-lin: Quadratic and linear-time MMD tests — L1-opt-ME --- L1-grid-ME









Sample test size

	L1-op-ME
<u> </u>	L1-grid-ME
	L1-opt-SCF
	L1-grid-SCF

	L1-opt-SCF
	L1-grid-SCF
• • • •	ME-full
	SCF-full
	MMD-lin

Higgs Dataset

