

Hypothesis Testing with Kernel Embeddings on Interdependent Data

Dino Sejdinovic

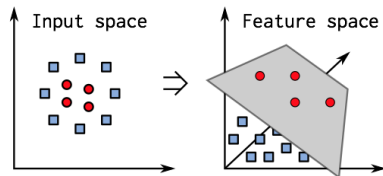
Department of Statistics
University of Oxford

joint work with
Kacper Chwialkowski and Arthur Gretton (Gatsby Unit, UCL)

9 April 2015
Dagstuhl

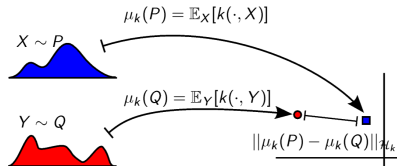
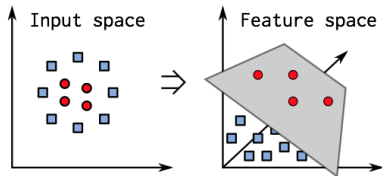
Kernel Embedding

- **feature map:** $x \mapsto k(\cdot, x) \in \mathcal{H}_k$
instead of
 $x \mapsto (\varphi_1(x), \dots, \varphi_s(x)) \in \mathbb{R}^s$
- $\langle k(\cdot, x), k(\cdot, y) \rangle_{\mathcal{H}_k} = k(x, y)$
inner products easily **computed**



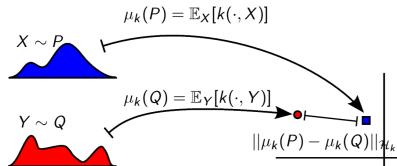
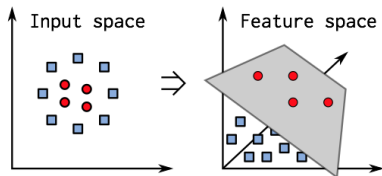
Kernel Embedding

- **feature map:** $x \mapsto k(\cdot, x) \in \mathcal{H}_k$
instead of
 $x \mapsto (\varphi_1(x), \dots, \varphi_s(x)) \in \mathbb{R}^s$
- $\langle k(\cdot, x), k(\cdot, y) \rangle_{\mathcal{H}_k} = k(x, y)$
inner products easily **computed**
- **embedding:**
 $P \mapsto \mu_k(P) = \mathbb{E}_{X \sim P} k(\cdot, X) \in \mathcal{H}_k$
instead of
 $P \mapsto (\mathbb{E}\varphi_1(X), \dots, \mathbb{E}\varphi_s(X)) \in \mathbb{R}^s$
- $\langle \mu_k(P), \mu_k(Q) \rangle_{\mathcal{H}_k} = \mathbb{E}_{X, Y} k(X, Y)$
inner products easily **estimated**



Kernel Embedding

- **feature map:** $x \mapsto k(\cdot, x) \in \mathcal{H}_k$
instead of
 $x \mapsto (\varphi_1(x), \dots, \varphi_s(x)) \in \mathbb{R}^s$
- $\langle k(\cdot, x), k(\cdot, y) \rangle_{\mathcal{H}_k} = k(x, y)$
inner products easily **computed**
- **embedding:**
 $P \mapsto \mu_k(P) = \mathbb{E}_{X \sim P} k(\cdot, X) \in \mathcal{H}_k$
instead of
 $P \mapsto (\mathbb{E} \varphi_1(X), \dots, \mathbb{E} \varphi_s(X)) \in \mathbb{R}^s$
- $\langle \mu_k(P), \mu_k(Q) \rangle_{\mathcal{H}_k} = \mathbb{E}_{X, Y} k(X, Y)$
inner products easily **estimated**
- $\mu_k(P)$ represents expectations w.r.t. P , i.e.,
 $\mathbb{E}_X f(X) = \mathbb{E}_X \langle f, k(\cdot, X) \rangle_{\mathcal{H}_k} = \langle f, \mu_k(P) \rangle_{\mathcal{H}_k} \quad \forall f \in \mathcal{H}_k$

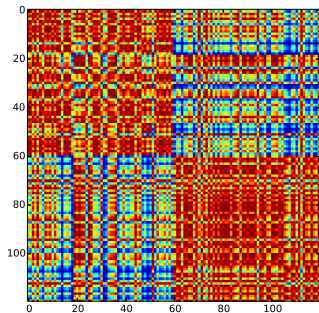
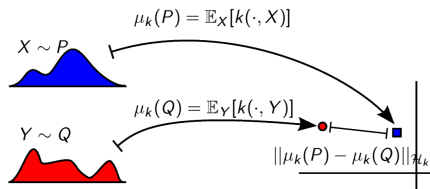


Kernel MMD

Definition

Kernel metric (MMD) between P and Q :

$$\begin{aligned}\text{MMD}_k(P, Q) &= \|\mathbb{E}_X k(\cdot, X) - \mathbb{E}_Y k(\cdot, Y)\|_{\mathcal{H}_k}^2 \\ &= \mathbb{E}_{XX'} k(X, X') + \mathbb{E}_{YY'} k(Y, Y') - 2\mathbb{E}_{XY} k(X, Y)\end{aligned}$$



Kernel MMD

- A polynomial kernel $k(x, x') = (1 + x^\top x')^s$ on \mathbb{R}^p captures the difference in first s (mixed) moments only
- For a certain family of kernels (**characteristic/universal**):
 $\text{MMD}_k(P, Q) = 0$ iff $P = Q$: Gaussian $\exp(-\frac{1}{2\sigma^2} \|z - z'\|_2^2)$,
Laplacian, inverse multiquadratics, B_{2n+1} -splines...
- Under mild assumptions, k -MMD metrizes weak* topology on probability measures (Sriperumbudur, 2010):

$$\text{MMD}_k(P_n, P) \rightarrow 0 \Leftrightarrow P_n \rightsquigarrow P$$

Nonparametric two-sample tests

- Testing $\mathbf{H}_0 : \mathbf{P} = \mathbf{Q}$ vs. $\mathbf{H}_A : \mathbf{P} \neq \mathbf{Q}$
based on samples $\{x_i\}_{i=1}^{n_x} \sim \mathbf{P}$, $\{y_i\}_{i=1}^{n_y} \sim \mathbf{Q}$.

- Test statistic is an estimate of

$$\text{MMD}_k(\mathbf{P}, \mathbf{Q}) = \mathbb{E}_{\mathbf{X}\mathbf{X}'} k(\mathbf{X}, \mathbf{X}') + \mathbb{E}_{\mathbf{Y}\mathbf{Y}'} k(\mathbf{Y}, \mathbf{Y}') - 2\mathbb{E}_{\mathbf{X}\mathbf{Y}} k(\mathbf{X}, \mathbf{Y}):$$

$$\widehat{\text{MMD}}_k = \frac{1}{n_x(n_x - 1)} \sum_{i \neq j} k(x_i, x_j) + \frac{1}{n_y(n_y - 1)} \sum_{i \neq j} k(y_i, y_j) - \frac{2}{n_x n_y} \sum_{i, j} k(x_i, y_j).$$

- Degenerate U-statistic: $\frac{1}{\sqrt{n}}$ -convergence to MMD under \mathbf{H}_A ,
 $\frac{1}{n}$ -convergence to 0 under \mathbf{H}_0 .
- $O(n^2)$ to compute ($n = n_x + n_y$) – various approximations (block-based, random features) trade computation for power.

Test threshold

- For i.i.d. data, under $\mathbf{H}_0 : \mathbf{P} = \mathbf{Q}$:

$$\frac{n_x n_y}{n_x + n_y} \widehat{\text{MMD}}_k \rightsquigarrow \sum_{r=1}^{\infty} \lambda_r (Z_r^2 - 1), \quad \{Z_r\}_{r=1}^{\infty} \stackrel{i.i.d.}{\sim} \mathcal{N}(0, 1)$$

- $\{\lambda_r\}$ depend on both k and \mathbf{P} : eigenvalues of $\mathbf{T} : L_2 \rightarrow L_2$,

$$(\mathbf{T}f)(x) \mapsto \int f(x') \underbrace{\tilde{k}(x, x')}_{\text{centred}} d\mathbf{P}(x').$$

- Asymptotic null distribution typically estimated using a permutation test.
- For interdependent samples, $\{Z_r\}_{r=1}^{\infty}$ are correlated, with the correlation structure dependent on the correlation structure within the samples.

Nonparametric independence tests

- $H_0 : X \perp\!\!\!\perp Y$
- $H_A : X \not\perp\!\!\!\perp Y$

Nonparametric independence tests

- $H_0 : X \perp\!\!\!\perp Y \Leftrightarrow \mathbf{P}_{XY} = \mathbf{P}_X \mathbf{P}_Y$
- $H_A : X \not\perp\!\!\!\perp Y \Leftrightarrow \mathbf{P}_{XY} \neq \mathbf{P}_X \mathbf{P}_Y$

- Test statistic:

$$\text{HSIC}(X, Y) = \left\| \mu_{\kappa}(\hat{P}_{XY}) - \mu_{\kappa}(\hat{P}_X \hat{P}_Y) \right\|_{\mathcal{H}_{\kappa}}^2,$$

with $\kappa = k \otimes l$

Gretton et al (2005, 2008); Smola et al (2007);

- Related to distance covariance (dCov) in statistics literature Szekely et al (AoS 2007, AoAS 2009); S. et al (AoS 2013)

$$k(\boxed{1}, \boxed{2}) \quad l(\boxed{1}, \boxed{2})$$

↓

$$\kappa(\boxed{1}, \boxed{1}, \boxed{2}, \boxed{2}) = k(\boxed{1}, \boxed{2}) \times l(\boxed{1}, \boxed{2})$$

HSIC computation

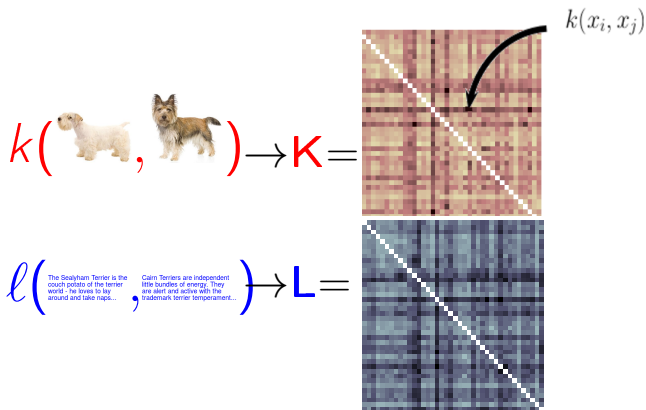
$$k(\text{img}_1, \text{img}_2)$$

$$l(\text{desc}_1, \text{desc}_2)$$

The Sealyham Terrier is the couch potato of the terrier world - he loves to lay around and take naps...

Calm Terriers are independent little bundles of energy. They are alert and active with the trademark terrier temperament...

HSIC computation

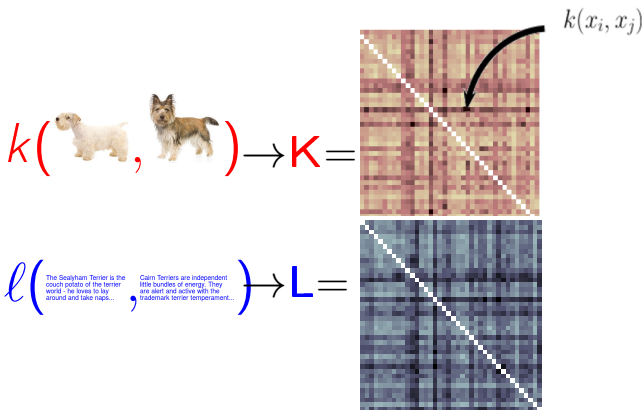


- **HSIC** measures *average similarity between the kernel matrices*:

$$\text{HSIC}(X, Y) = \frac{1}{n^2} \langle H\mathbf{K}H, H\mathbf{L}H \rangle$$

- $H = I - \frac{1}{n}\mathbf{1}\mathbf{1}^\top$
(centering matrix)

HSIC computation



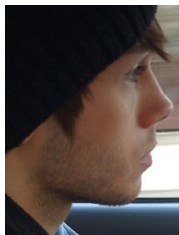
- **HSIC** measures *average similarity between the kernel matrices*:

$$\text{HSIC}(X, Y) = \frac{1}{n^2} \langle H\mathbf{K}H, H\mathbf{L}H \rangle$$

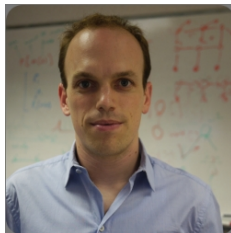
- $H = I - \frac{1}{n}\mathbf{1}\mathbf{1}^\top$
(centering matrix)

Extensions: conditional independence testing (Fukumizu, Gretton, Sun and Schölkopf, 2008; Zhang, Peters, Janzing and Schölkopf, 2011), three-variable interaction / V-structure discovery (S., Gretton and Bergsma, 2013)

Kernel tests on time series

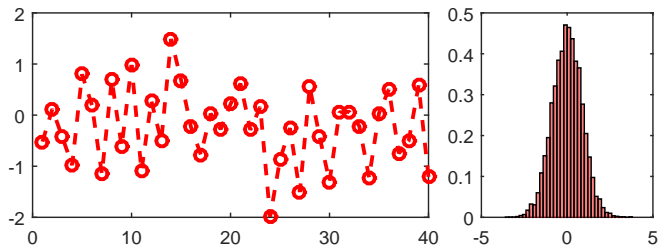
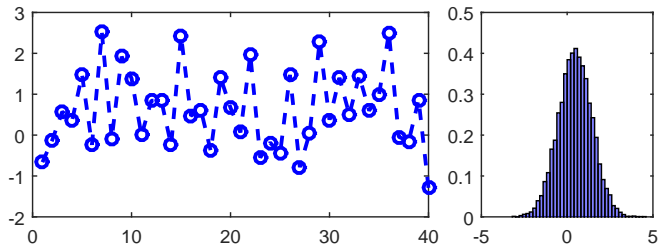


Kacper Chwialkowski



Arthur Gretton

Test calibration for dependent observations



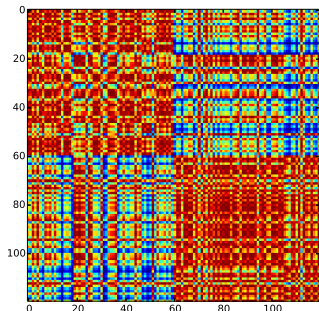
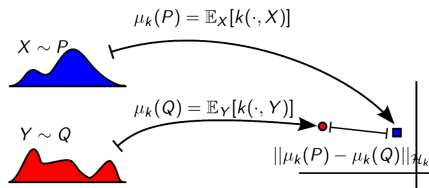
Is
 P
the same
distribution as
 Q
?

Kernel MMD

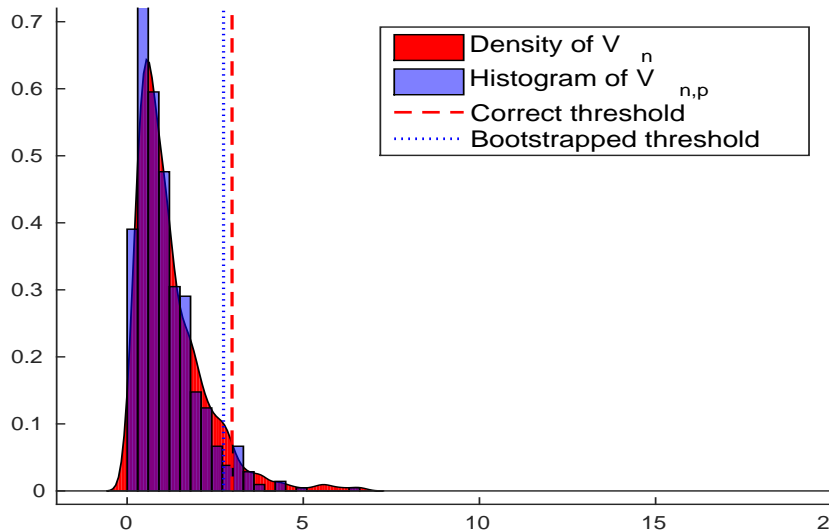
Definition

Kernel metric (MMD) between P and Q :

$$\begin{aligned} \text{MMD}_k(P, Q) &= \|\mathbb{E}_X k(\cdot, X) - \mathbb{E}_Y k(\cdot, Y)\|_{\mathcal{H}_k}^2 \\ &= \mathbb{E}_{XX'} k(X, X') + \mathbb{E}_{YY'} k(Y, Y') - 2\mathbb{E}_{XY} k(X, Y) \end{aligned}$$

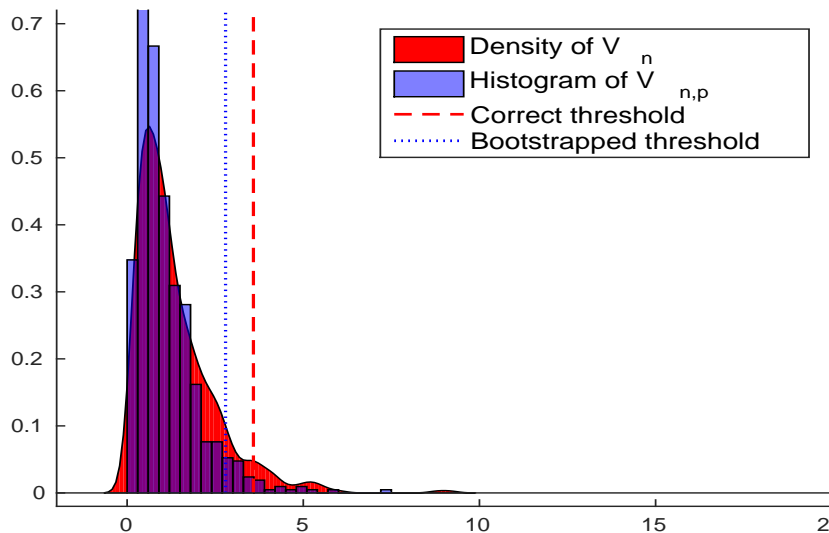


Permutation test on AR(1): $X_{t+1} = aX_t + \sqrt{1 - a^2}\epsilon_t$



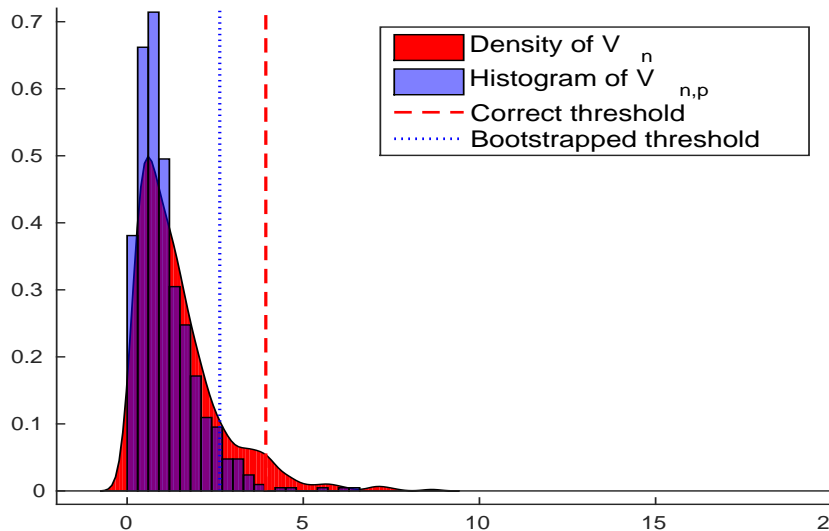
$a = 0.1$

Permutation test on AR(1): $X_{t+1} = aX_t + \sqrt{1 - a^2}\epsilon_t$



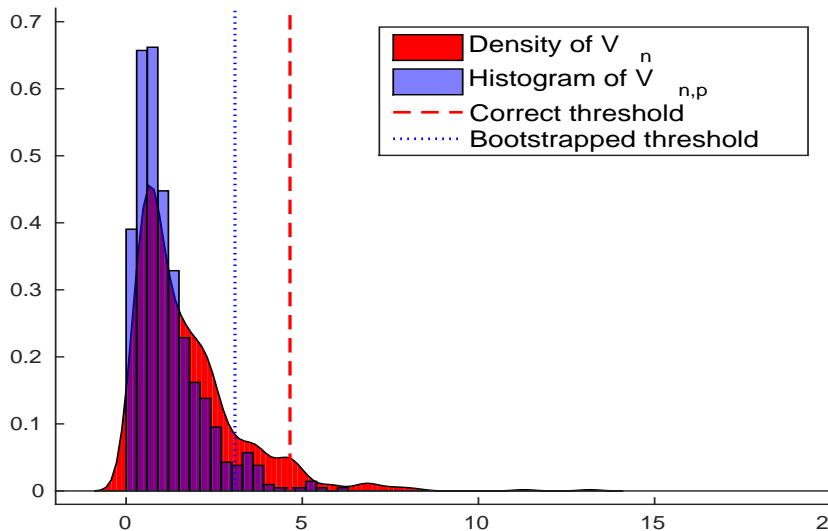
$a = 0.2$

Permutation test on AR(1): $X_{t+1} = aX_t + \sqrt{1 - a^2}\epsilon_t$



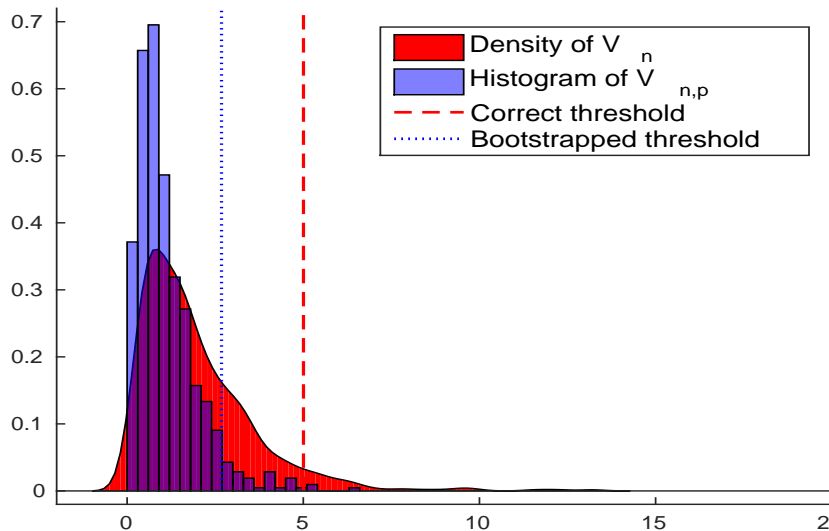
$a = 0.3$

Permutation test on AR(1): $X_{t+1} = aX_t + \sqrt{1 - a^2}\epsilon_t$



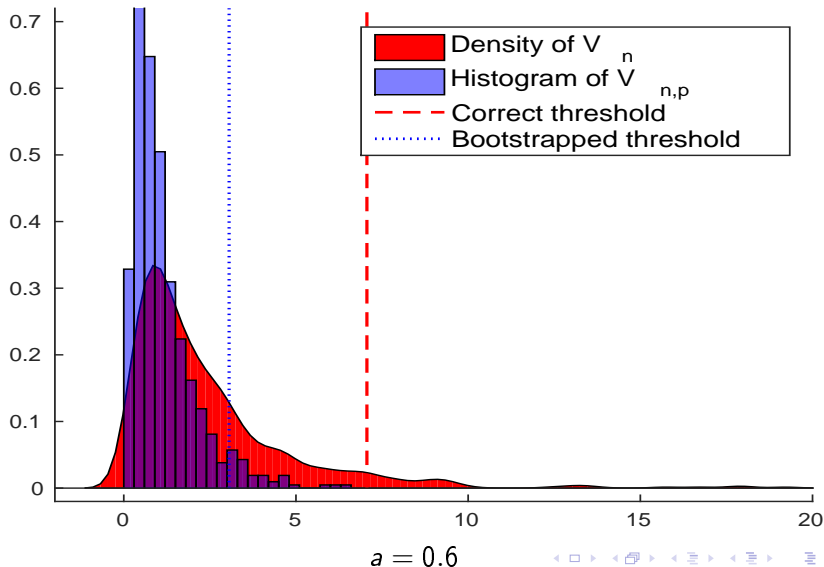
$a = 0.4$

Permutation test on AR(1): $X_{t+1} = aX_t + \sqrt{1 - a^2}\epsilon_t$

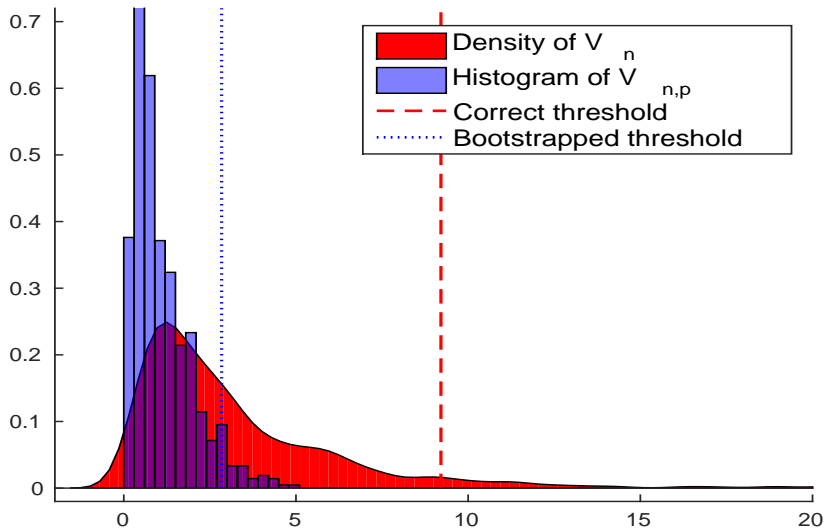


$a = 0.5$

Permutation test on AR(1): $X_{t+1} = aX_t + \sqrt{1 - a^2}\epsilon_t$

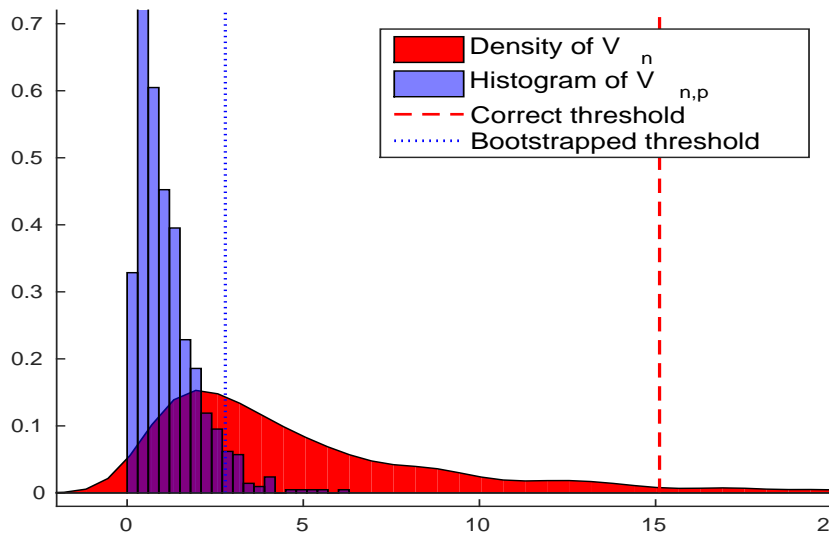


Permutation test on AR(1): $X_{t+1} = aX_t + \sqrt{1 - a^2}\epsilon_t$



$a = 0.7$

Permutation test on AR(1): $X_{t+1} = aX_t + \sqrt{1 - a^2}\epsilon_t$



$a = 0.8$

Wild Bootstrap

Wild bootstrap process (Leucht and Neumann, 2013):

$W_{t,n} = e^{-1/l_n} W_{t-1,n} + \sqrt{1 - e^{-2/l_n}} \epsilon_t$ where $W_{0,n}, \epsilon_1, \dots, \epsilon_n \stackrel{i.i.d.}{\sim} \mathcal{N}(0, 1)$, and $\tilde{W}_{t,n} = W_{t,n} - \frac{1}{n} \sum_{j=1}^n W_{j,n}$.

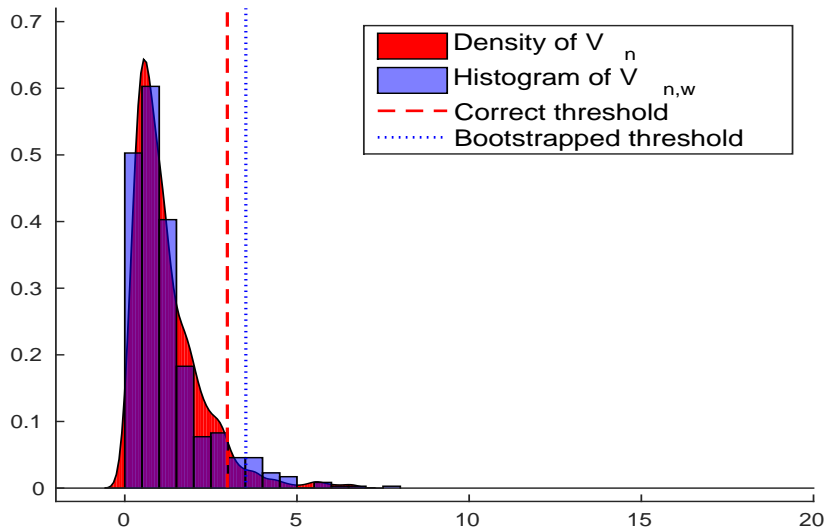
$$\widehat{\text{MMD}}_{k,wb} := \frac{1}{n_x^2} \sum_{i=1}^{n_x} \sum_{j=1}^{n_x} \tilde{W}_{i,n_x}^{(x)} \tilde{W}_{j,n_x}^{(x)} k(x_i, x_j) - \frac{1}{n_x^2} \sum_{i=1}^{n_y} \sum_{j=1}^{n_y} \tilde{W}_{i,n_y}^{(y)} \tilde{W}_{j,n_y}^{(y)} k(y_i, y_j) \\ - \frac{2}{n_x n_y} \sum_{i=1}^{n_x} \sum_{j=1}^{n_y} \tilde{W}_{i,n_x}^{(x)} \tilde{W}_{j,n_y}^{(y)} k(x_i, y_j).$$

Theorem (Chwialkowski, S. and Gretton, 2014)

Let k be bounded and Lipschitz continuous, and let $\{X_t\} \sim P$ and $\{Y_t\} \sim Q$ both be τ -dependent with $\tau(r) = O(r^{-6-\epsilon})$, but independent of each other.

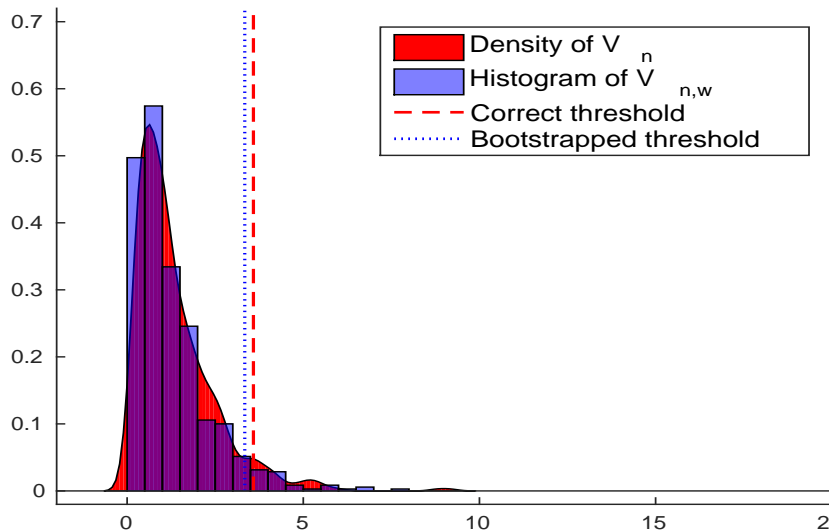
Then, under $\mathbf{H}_0 : P = Q$, $\varphi \left(\frac{n_x n_y}{n_x + n_y} \widehat{\text{MMD}}_k, \frac{n_x n_y}{n_x + n_y} \widehat{\text{MMD}}_{k,b} \right) \xrightarrow{P} 0$ as $n_x, n_y \rightarrow \infty$, where φ is the Prokhorov metric.

Wild Bootstrap on AR(1): $X_{t+1} = aX_t + \sqrt{1 - a^2}\epsilon_t$



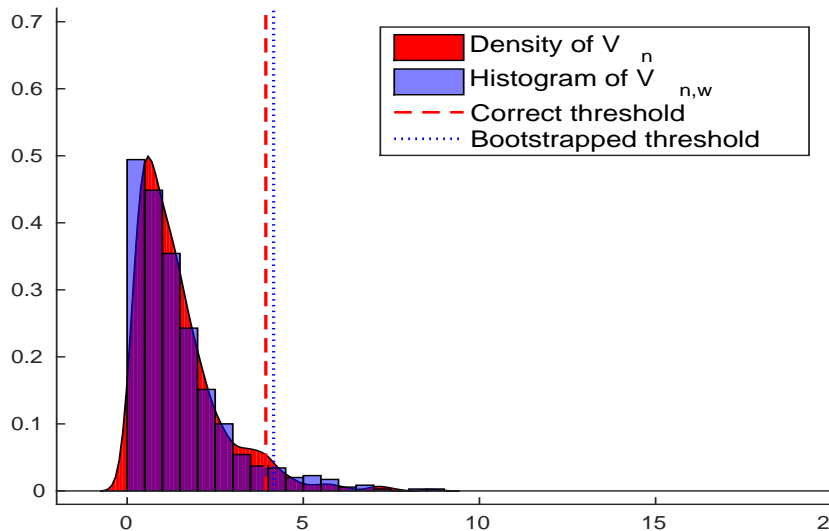
$a = 0.1$

Wild Bootstrap on AR(1): $X_{t+1} = aX_t + \sqrt{1 - a^2}\epsilon_t$



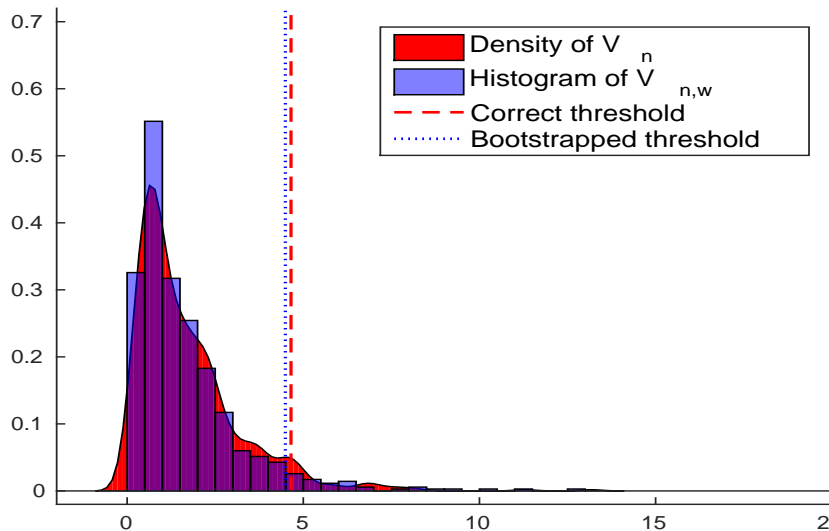
$a = 0.2$

Wild Bootstrap on AR(1): $X_{t+1} = aX_t + \sqrt{1 - a^2}\epsilon_t$



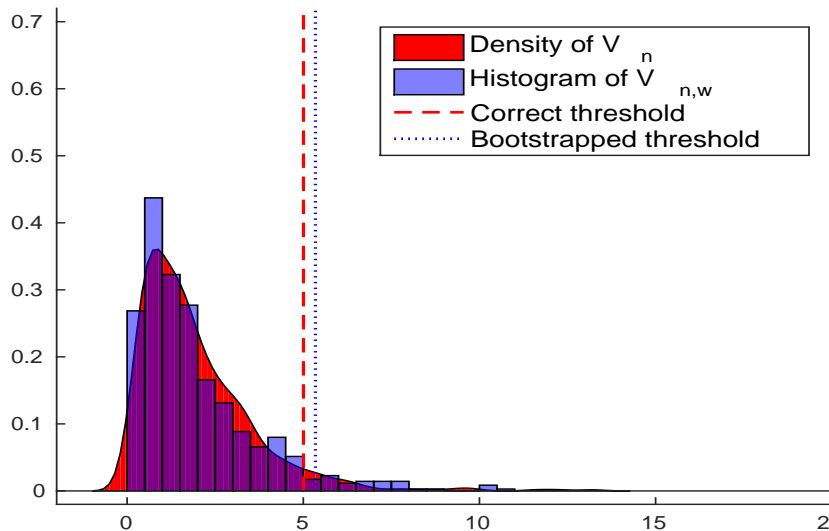
$a = 0.3$

Wild Bootstrap on AR(1): $X_{t+1} = aX_t + \sqrt{1 - a^2}\epsilon_t$



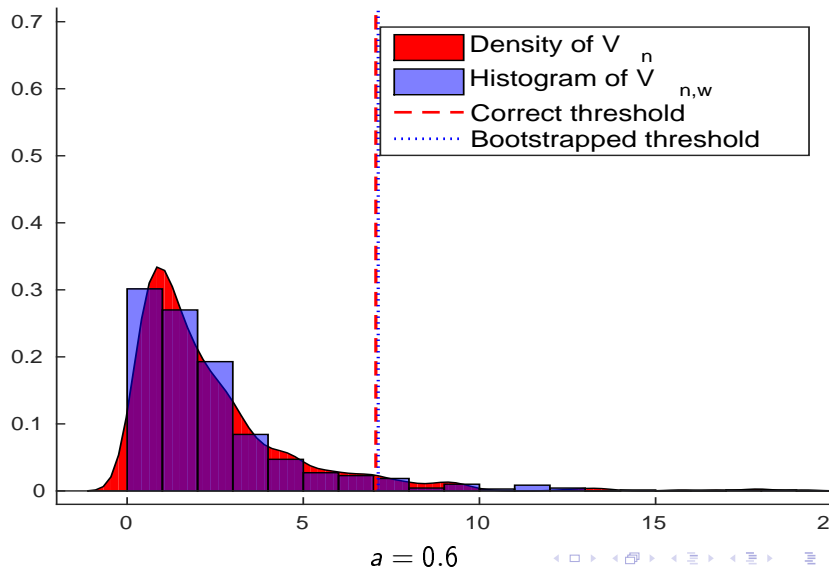
$a = 0.4$

Wild Bootstrap on AR(1): $X_{t+1} = aX_t + \sqrt{1 - a^2}\epsilon_t$

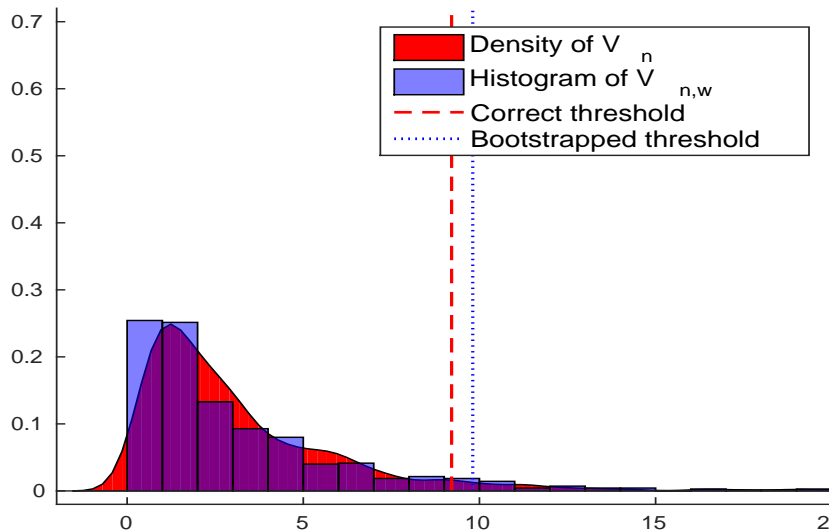


$a = 0.5$

Wild Bootstrap on AR(1): $X_{t+1} = aX_t + \sqrt{1 - a^2}\epsilon_t$

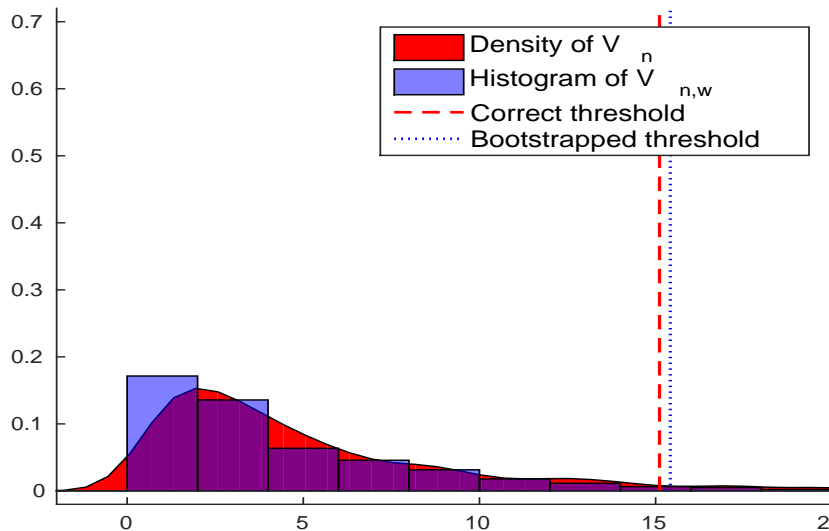


Wild Bootstrap on AR(1): $X_{t+1} = aX_t + \sqrt{1 - a^2}\epsilon_t$



$a = 0.7$

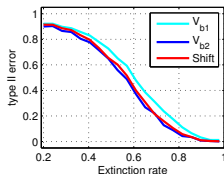
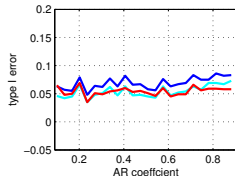
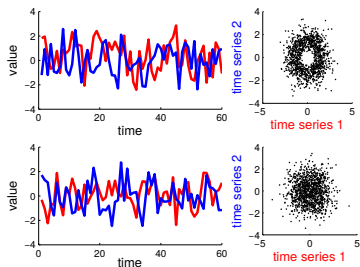
Wild Bootstrap on AR(1): $X_{t+1} = aX_t + \sqrt{1 - a^2}\epsilon_t$



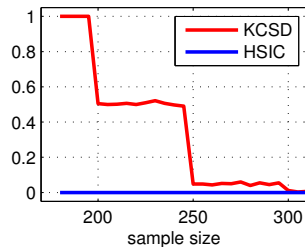
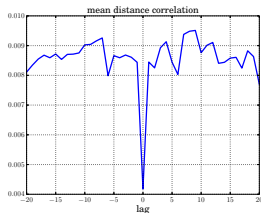
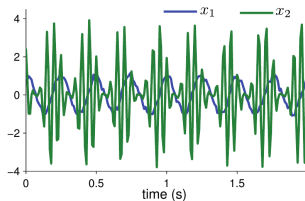
$a = 0.8$

Test calibration for dependent observations

Two-sample test	experiment \ method	perm.	wild
MCMC diagnostics	i.i.d. vs i.i.d. (H_0)	.040	.012
	i.i.d. vs Gibbs (H_0)	.528	.052
	Gibbs vs Gibbs (H_0)	.680	.060



Time Series Coupled at a Lag



$$X_t = \cos(\phi_{t,1}), \quad \phi_{t,1} = \phi_{t-1,1} + 0.1\epsilon_{1,t} + 2\pi f_1 T_s, \quad \epsilon_{1,t} \stackrel{i.i.d.}{\sim} \mathcal{N}(0, 1),$$
$$Y_t = [2 + C \sin(\phi_{t,1})] \cos(\phi_{t,2}), \quad \phi_{t,2} = \phi_{t-1,2} + 0.1\epsilon_{2,t} + 2\pi f_2 T_s, \quad \epsilon_{2,t} \stackrel{i.i.d.}{\sim} \mathcal{N}(0, 1).$$

Parameters: $C = 0.4$, $f_1 = 4\text{Hz}$, $f_2 = 20\text{Hz}$, $\frac{1}{T_s} = 100\text{Hz}$.

- M. Besserve, N.K. Logothetis, and B. Schölkopf. **Statistical analysis of coupled time series with kernel cross-spectral density operators**. *NIPS 2013*.

Summary

- Interdependent data lead to incorrect Type I control for kernel tests (too many false positives).
- Consistency of a wild bootstrap procedure under weak long-range dependencies (τ -mixing), applicable to both two-sample and independence tests
- Applications: MCMC diagnostics, time series dependence across multiple lags

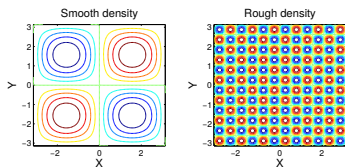
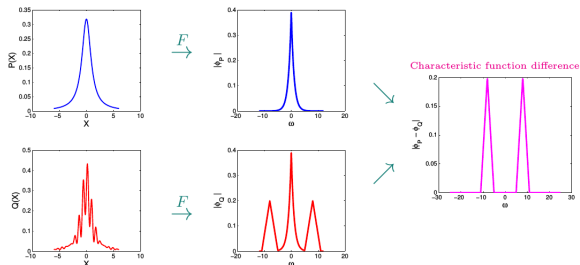
Open questions

- Interdependent case: how to select parameters of the wild bootstrap / block bootstrap - requires estimating mixing properties of the time series first?
- Large-scale testing: tradeoffs between computation and power
- How to interpret the discovered differences in distributions / discovered dependence? Do we really care about all possible differences between distributions?
- Tuning parameters - can select kernels/hyperparameters to directly optimize relative efficiency of the test, but how does this affect tradeoffs with interdependent data? Sensitive interplay between the kernel hyperparameter and the wild bootstrap parameters
- Multivariate interaction and graphical model selection - approximations?

References

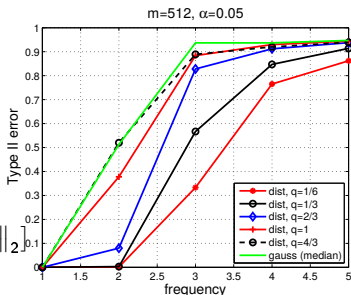
- K. Chwialkowski, DS and A. Gretton, **A wild bootstrap for degenerate kernel tests**. *Advances in Neural Information Processing Systems (NIPS)* 27, Dec. 2014.
- DS, B. Sriperumbudur, A. Gretton and K. Fukumizu, **Equivalence of distance-based and RKHS-based statistics in hypothesis testing**. *Ann. Statist.* 41(5): 2263-2291, Oct. 2013.
- M. Besserve, N.K. Logothetis and B. Schölkopf, **Statistical analysis of coupled time series with kernel cross-spectral density operators**. *Advances in Neural Information Processing Systems (NIPS)* 26, Dec. 2013.
- A. Leucht and M.H. Neumann, **Dependent wild bootstrap for degenerate U- and V-statistics**. *J. Multivar. Anal.* 117:257-280, 2013.
- A. Gretton, K.M. Borgwardt, M.J. Rasch, B. Schölkopf and A. Smola, **A Kernel Two-Sample Test**. *J. Mach. Learn. Res.* 13(Mar): 723-773, 2012.

Kernels and characteristic functions



E-distance/dCov of Székely and Rizzo (2004,2005), Székely et al (2007):

$$\begin{aligned} \nu^2(X, Y) &= \mathbb{E}_{X,Y} \mathbb{E}_{X',Y'} \|X - X'\|_2 \|Y - Y'\|_2 \\ &\quad + \mathbb{E}_X \mathbb{E}_{X'} \|X - X'\|_2 \mathbb{E}_Y \mathbb{E}_{Y'} \|Y - Y'\|_2 \\ &\quad - 2 \mathbb{E}_{X,Y} [\mathbb{E}_{X'} \|X - X'\|_2 \mathbb{E}_{Y'} \|Y - Y'\|_2] \end{aligned}$$



DS, B. Sriperumbudur, A. Gretton and K. Fukumizu, **Equivalence of distance-based and RKHS-based statistics in hypothesis testing**. *Annals of Statistics* 41(5), p. 2263-2291, 2013.

Embeddings in Mercer's Expansion

Mercer's Expansion

For a compact metric space \mathcal{X} , and a continuous kernel k ,

$$k(x, y) = \sum_{r=1}^{\infty} \lambda_r \Phi_r(x) \Phi_r(y),$$

with $\{\lambda_r, \Phi_r\}_{r \geq 1}$ eigenvalue, eigenfunction pairs of $f \mapsto \int f(x)k(\cdot, x)dP(x)$ on $L_2(P)$.

$$\mathcal{H}_k \ni k(\cdot, x) \leftrightarrow \{\sqrt{\lambda_r} \Phi_r(x)\} \in \ell_2$$

$$\mathcal{H}_k \ni \mu_k(P) \leftrightarrow \{\sqrt{\lambda_r} \mathbb{E} \Phi_r(X)\} \in \ell_2$$

$$\left\| \mu_k(\hat{P}) - \mu_k(\hat{Q}) \right\|_{\mathcal{H}_k}^2 = \sum_{r=1}^{\infty} \lambda_r \left(\frac{1}{n_x} \sum_{t=1}^{n_x} \Phi_r(X_t) - \frac{1}{n_y} \sum_{t=1}^{n_y} \Phi_r(Y_t) \right)^2$$

Wild Bootstrap

- $\rho_x = n_x/n$, $\rho_y = n_y/n$
- $\{W_{t,n}\}_{1 \leq t \leq n}$, $\mathbb{E}W_{t,n} = 0$, $\mathbb{E}[W_{t,n}W_{t',n}] = \zeta\left(\frac{|t'-t|}{\ell_n}\right)$, with $\lim_{u \rightarrow 0} \zeta(u) \rightarrow 1$

$$\rho_x \rho_y n \widehat{\text{MMD}}_k = \sum_{r=1}^{\infty} \lambda_r \left(\sqrt{\rho_y} \sum_{t=1}^{n_x} \frac{\Phi_r(X_t)}{\sqrt{n_x}} - \sqrt{\rho_x} \sum_{t=1}^{n_y} \frac{\Phi_r(Y_t)}{\sqrt{n_y}} \right)^2$$

$$\rho_x \rho_y n \widehat{\text{MMD}}_{k,wb} = \sum_{r=1}^{\infty} \lambda_r \left(\sqrt{\rho_y} \sum_{t=1}^{n_x} \frac{\Phi_r(X_t) \tilde{W}_{t,n_x}^{(y)}}{\sqrt{n_x}} - \sqrt{\rho_x} \sum_{t=1}^{n_y} \frac{\Phi_r(Y_t) \tilde{W}_{t,n_y}^{(y)}}{\sqrt{n_y}} \right)^2$$

- $\mathbb{E}[\Phi_r(X_1)W_{1,n}\Phi_s(X_t)W_{t,n}] = \mathbb{E}[\Phi_r(X_1)\Phi_s(X_t)]\zeta\left(\frac{|t-1|}{\ell_n}\right) \xrightarrow{n \rightarrow \infty} \mathbb{E}[\Phi_r(X_1)\Phi_s(X_t)]$, $\forall t, r, s$ provided dependence between X_1 and X_t “disappears fast enough” (a τ -mixing condition).

ICML Workshop on Large-Scale Kernel Learning

Lille, France, 11 July 2015 (collocated with ICML 2015)

- Foundational algorithmic techniques for large-scale kernel learning: matrix factorization, randomization and approximation, variational inference and sampling, inducing variables, random Fourier features, unifying frameworks
- Interface between kernel methods and deep learning architectures
- Tradeoffs between statistical and computational efficiency in kernel methods
- Stochastic gradient techniques with kernel methods
- Large-scale multiple kernel learning
- Large-scale representation learning with kernels
- Large-scale kernel methods for complex data types beyond perceptual data
- Confirmed speakers: Francis Bach, Neil Lawrence, Russ Salakhutdinov, Marius Kloft, Zaid Harchaoui
- Deadline for Submissions: **Friday, May 1st, 2015, 23:00 UTC.**