

Kernel Methods for Computing Probability

Distributions

Problem: testing for differences and/or dependences btw r.v.'s in high dimensional spaces ("curse of dimensionality") with limited samples is very difficult

By working in ∞ -dimensional spaces, kernel methods allow us to do this relatively well for arbitrary data structures (e.g. discrete/continuous)

Key: Mean embedding:

"kernel trick": $f(x) = \langle f, \phi_x \rangle_{\mathcal{F}}$

\Rightarrow "mean trick": $\mathbb{E}_{x \sim P} [f(x)] = \langle \mu_P, f \rangle_{\mathcal{F}}$

Empirically, you can estimate μ_P via:

$$\hat{\mu}_P = \frac{1}{N} \sum_{i=1}^N \phi_{x_i}, \quad x_i \stackrel{i.i.d.}{\sim} P$$

you can prove μ_P exists via the Riesz theorem

More generally,
 $\mu_P(x) = \langle \mu_P, \phi_x \rangle_{\mathcal{F}}$
 $= \langle \mu_P, K(\cdot, x) \rangle_{\mathcal{F}}$
 by definition ("mean trick") $= \mathbb{E}_{x \sim P} K(x, x)$
 $= K(x, \mathbb{E}_{x \sim P} x)$

Note that $\hat{\mu}_P$ is an infinitely dimensional vector in the RKHS \mathcal{F} , i.e. it is a function! So, we can get estimates of the mean in different regions of \mathcal{X} by indexing via the function argument.

Analogously to how we do this in finite space:

For $\{x^{(i)}\} \subseteq \mathbb{R}^d$, $\bar{x}_d = \frac{1}{N} \sum_i x_d^{(i)} = \frac{1}{N} \sum_i x^{(i)T} e_d$

↑ find mean on dth dimension ↑ projection onto dth basis vector

$\begin{bmatrix} 0 \\ \vdots \\ 1 \\ \vdots \\ 0 \end{bmatrix} \begin{matrix} \vdots \\ d \\ \vdots \\ 0 \end{matrix}$

probability feature map

For $\{x^{(i)}\} \subseteq \mathcal{X}$, $\mu_P(d) = \frac{1}{N} \sum_i \langle \phi_{x^{(i)}}, \phi_d \rangle = \frac{1}{N} \sum_i K(x^{(i)}, d)$

$d \in \mathcal{X}$

expectation over a kernel!

We can formally compare means in feature space via the Maximum mean discrepancy:

$$\text{MMD}(P, Q; \mathcal{F}) = \sup_{f \in \mathcal{F}} [\mathbb{E}_P f(X) - \mathbb{E}_Q f(Y)], \quad X \sim P, Y \sim Q$$

→ MMD = 0 iff $P=Q$ when \mathcal{F} is the unit ball in a characteristic RKHS

$$= \sup_{f \in \mathcal{F}} [\langle f, \mu_P \rangle_{\mathcal{F}} - \langle f, \mu_Q \rangle_{\mathcal{F}}]$$

$$= \sup_{f \in \mathcal{F}} \langle f, \mu_P - \mu_Q \rangle_{\mathcal{F}}$$

$$= \|\mu_P - \mu_Q\|_{\mathcal{F}}$$

But, of course, we never have access to P and Q , only samples. So, we estimate MMD as:

$$\text{MMD}^2 = \|\mu_P - \mu_Q\|^2 = \langle \mu_P - \mu_Q, \mu_P - \mu_Q \rangle$$

$$= \langle \mu_P, \mu_P \rangle + \langle \mu_Q, \mu_Q \rangle - 2 \langle \mu_P, \mu_Q \rangle$$

$$\approx \frac{1}{N(N-1)} \sum_{i=1}^N \sum_{j \neq i}^N k(x_i, x_j) + \frac{1}{M(M-1)} \sum_{i=1}^M \sum_{j \neq i}^M k(y_i, y_j) - \frac{2}{NM} \sum_{i,j} k(x_i, y_j)$$

where $x_i \stackrel{i.i.d.}{\sim} P, y_i \stackrel{i.i.d.}{\sim} Q$

We now want to make sure that MMD does what we want: i.e. $\text{MMD} = 0$ iff $P=Q$. As mentioned above, this holds whenever \mathcal{F} is a characteristic RKHS

For a translation invariant ~~kernel~~ and periodic kernel k (see below) that ~~is~~ k is a characteristic kernel (with associated characteristic RKHS \mathcal{F}) if $\forall \mu \hat{\mu} \neq 0$:

Supposing $p(x)$ is defined only over $x \in [-\pi, \pi]$,

$$\mu_P(z) = \mathbb{E}_{x \sim P} k(z, x)$$

$$= \mathbb{E}_{x \sim P} k(z-x)$$

$$= \int_{-\pi}^{\pi} k(z-x) p(x) dx$$

k is translation invariant

k is periodic

The Fourier transform of μ_p is then given by:

$$\hat{\mu}_{p,l} = \int_{-\pi}^{\pi} k(z-x) e^{-ilz} dP(x) dz$$

Let $v = z-x \Leftrightarrow dv = dz$

$$= \int_{-\pi}^{\pi} k(v) e^{-il(x+v)} dP(x) dv$$

$$= \int k(v) e^{-ilv} dv \int_{-\pi}^{\pi} e^{-ilx} dP(x)$$

$$= \hat{K}_l \psi_{p,l}$$

Recall now that for an RKHS \mathcal{F} associated with a periodic and translation invariant kernel k , $\|f\|_{\mathcal{F}}^2 = \sum_{l=-\infty}^{\infty} \frac{|\hat{f}_l|^2}{K_l}$

So, we now have

$$\text{MMD}^2 = \|\mu_p - \mu_q\|_{\mathcal{F}}^2 = \sum_{l=-\infty}^{\infty} \frac{|\hat{K}_l \psi_{p,l} - \hat{K}_l \psi_{q,l}|^2}{\hat{K}_l}$$

$$= \sum_{l=-\infty}^{\infty} |\psi_{p,l} - \psi_{q,l}|^2 \hat{K}_l$$

Two points:

① Note that as long as $\hat{K}_l \neq 0$ for all l , $\text{MMD}^2 = 0$ only when $p = q$, making k a characteristic kernel

② Remember that \hat{K}_l decays with l , so this definition of MMD can be interpreted as penalizing differences in lower frequencies more than in higher frequencies

~~For an arbitrary kernel $k(x,y)$ and probability distribution $P(x)$ on \mathbb{R}^D , we need a similar proof by invoking~~

For an arbitrary translation invariant kernel $k(x,y)$ on \mathbb{R}^D ,

we invoke Bochner's theorem: $k(x,y) = k(x-y) = k(z) = \int_{\mathbb{R}^D} e^{-iz \cdot \omega} d\Delta(\omega)$
 $\Rightarrow k$ characteristic for prob. measure on \mathbb{R}^D iff $\text{supp}(\Delta) = \mathbb{R}^D$
 Fourier transform of k

Having established MMD as a good measure of similarity b/w probability distributions, lets now think about how we can use it to ~~make~~ make inferences from data via hypothesis testing.

$$H_0: P = Q$$

$$H_1: P \neq Q$$

given $\{x_i\}_{i=1}^N, x_i \stackrel{i.i.d.}{\sim} P$

$\{y_j\}_{j=1}^M, y_j \stackrel{i.i.d.}{\sim} Q$

We want to reject H_0 if MMD is far from zero. We can then define this by picking a threshold criterion based on the asymptotic distribution of \widehat{MMD}^2 ($N \rightarrow \infty$)

$$\widehat{MMD}^2 = \frac{1}{N(N-1)} \sum_{i \neq j} \underbrace{k(x_i, x_j) - k(x_i, y_j) - k(y_i, x_j) + k(y_i, y_j)}_{h((x_i, y_i), (x_j, y_j))}$$

for $N \geq M, N \rightarrow \infty$

* note that \widehat{MMD} is missing some terms from our previous unbiased empirical estimate of MMD: \widehat{MMD} is still an unbiased estimator, but no longer minimum variance

When $P \neq Q$,

$$\sqrt{N} (\widehat{MMD}^2 - MMD^2) \sim \mathcal{N}(0, \sigma_u^2)$$

$$\text{where } \sigma_u^2 = 4 \left(\mathbb{E}_{z, z'} \left[\mathbb{E}_{z'} \left[h(z, z') \right]^2 \right] - \left[\mathbb{E}_{z, z'} h(z, z') \right]^2 \right)$$

$z = (x_i, y_i)$

When $P = Q$,

$$N \cdot \widehat{MMD}^2 \sim \sum_{l=1}^{\infty} \lambda_l (z_l^2 - 2)$$

this is an infinite sum of z^2 distributions

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

$$\int_{\mathcal{X}} \tilde{K}(x, x') \phi_l(x) dP(x) = \lambda_l \phi_l(x')$$

\mathcal{X} (centered)

\Rightarrow this is a degenerate U-statistic, can't compute the null distribution!

Pearson moment matching

permutation methods

Can we use a similar kernel-based metric to test for dependence b/w two random variables?

If X, Y are independent, then

$$P(X, Y) = P(X) P(Y) \quad (*)$$

Let's call these three distributions $P_{X,Y}, P_X, P_Y$, which are probability measures on $X \times Y, X, Y$ respectively.

We now move into arbitrary feature space by assuming two RKHS's \mathcal{F} and \mathcal{G} w/ associated kernels k and l such that $\mu_{P_X} \in \mathcal{F}, \mu_{P_Y} \in \mathcal{G}$.

If $X \sim P_X$ and $Y \sim P_Y$ are independent, then, for a function $f(x, y)$ with $x \in X$ and $y \in Y$ eqn. (*) tells us that

$$\mathbb{E}_{P_{X,Y}} f = \mathbb{E}_{P_X P_Y} f, \text{ where } f \text{ is a member of the Hilbert space } \mathcal{F} \times \mathcal{G}$$

$$\Leftrightarrow \mathbb{E}_{P_{X,Y}} f - \mathbb{E}_{P_X P_Y} f = 0 \quad (**)$$

We can then use an analogy of MMD ~~metric~~ called the Hilbert-Schmidt Independence Criterion (HSIC):

$$\begin{aligned} \text{HSIC}(P_{X,Y}, P_X P_Y) &= \text{MMD}(P_{X,Y}, P_X P_Y; \mathcal{F} \times \mathcal{G})^2 \\ &= \left(\sup_{\|f\|=1} \mathbb{E}_{P_{X,Y}} f - \mathbb{E}_{P_X P_Y} f \right)^2 \\ &= \| \mu_{P_{X,Y}} - \mu_{P_X P_Y} \|_{\mathcal{F} \times \mathcal{G}}^2 \end{aligned}$$

which should equal 0 if X, Y are independent (eqn. (**)).

Noting that

$$\begin{aligned} \mu_{P_{X,Y}}(\phi \otimes \psi) &= \mathbb{E}_{P_{X,Y}}(\phi(x)\psi(y)) = \mathbb{E}_{P_X P_Y}(\phi(x)\psi(y)) = \mu_{P_X P_Y}(\phi \otimes \psi) \\ \mu_{P_X P_Y}(\phi \otimes \psi) &= \mathbb{E}_{P_X P_Y}(\phi(x)\psi(y)) = \mathbb{E}_{P_X}(\phi(x)) \mathbb{E}_{P_Y}(\psi(y)) \end{aligned}$$

(Note: The above equations are crossed out in the original image with scribbles. The scribbles include the text "same kernel" and "Hilbert space".)

We can expand the HSIC into:

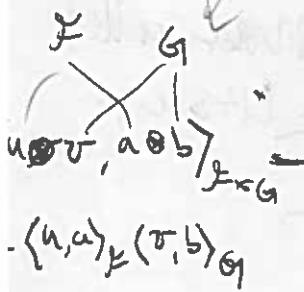
$$\begin{aligned}
 \text{HSIC}(P_{XY}, P_X P_Y) &= \|M_{P_{XY}} - M_{P_X P_Y}\|^2 \\
 &= \langle M_{P_{XY}}, M_{P_{XY}} \rangle + \langle M_{P_X P_Y}, M_{P_X P_Y} \rangle - 2 \langle M_{P_{XY}}, M_{P_X P_Y} \rangle \\
 &= \mathbb{E}_{x, y \sim P_{XY}} \mathbb{E}_{x', y' \sim P_{XY}} K(x, x') l(y, y') + \mathbb{E}_{x \sim P_X} \mathbb{E}_{x' \sim P_X} K(x, x') \mathbb{E}_{y \sim P_Y} \mathbb{E}_{y' \sim P_Y} l(y, y') \\
 &\quad - 2 \mathbb{E}_{x, y \sim P_{XY}} \left[\mathbb{E}_{x' \sim P_X} K(x, x') \right] \left[\mathbb{E}_{y' \sim P_Y} l(y, y') \right]
 \end{aligned}$$

For scratched out part on other side:

$$M_{P_{XY}} := \mathbb{E}_{x, y \sim P_{XY}} \phi_x \otimes \psi_y$$

$$\begin{aligned}
 M_{P_{XY}}(x, y) &= \langle M_{P_{XY}}, \phi_x \otimes \psi_y \rangle_{\mathcal{F} \otimes \mathcal{G}} \\
 &= \langle \mathbb{E}_{x', y' \sim P_{XY}} \phi_{x'} \otimes \psi_{y'}, \phi_x \otimes \psi_y \rangle_{\mathcal{F} \otimes \mathcal{G}} \\
 &= \mathbb{E}_{x', y' \sim P_{XY}} \langle \phi_{x'} \otimes \psi_{y'}, \phi_x \otimes \psi_y \rangle_{\mathcal{F} \otimes \mathcal{G}} \\
 &\stackrel{\text{linearity}}{=} \mathbb{E}_{x', y' \sim P_{XY}} \langle \phi_{x'}, \phi_x \rangle_{\mathcal{F}} \langle \psi_{y'}, \psi_y \rangle_{\mathcal{G}} \\
 &= \mathbb{E}_{x', y' \sim P_{XY}} K(x, x') l(y, y')
 \end{aligned}$$

you can show that $M_{P_{XY}}$ exists via Riesz theorem: exists in a Hilbert space $\mathcal{H}(\mathcal{F}, \mathcal{G})$ w/ Hilbert-Schmidt norm, not product. ↳ from which you prove



$$M_{P_X P_Y} := \mathbb{E}_{x \sim P_X} \mathbb{E}_{y \sim P_Y} \phi_x \otimes \psi_y$$

As above,

$$\begin{aligned}
 M_{P_X P_Y}(x, y) &= \mathbb{E}_{x' \sim P_X} \mathbb{E}_{y' \sim P_Y} \langle \phi_{x'} \otimes \psi_{y'}, \phi_x \otimes \psi_y \rangle \\
 &= \mathbb{E}_{x' \sim P_X} K(x, x') \mathbb{E}_{y' \sim P_Y} l(y, y')
 \end{aligned}$$

From this, you can show the above expansion of the HSIC.

Given a set of data points $\{(x_i, y_i)\}_{i=1}^N$, we can estimate the HSIC with

$$\hat{M}_{P_X} = \frac{1}{N} \sum_{i=1}^N \phi_{x_i}, \quad \hat{M}_{P_Y} = \frac{1}{N} \sum_{i=1}^N \psi_{y_i}, \quad \hat{M}_{P_{XY}} = \frac{1}{N} \sum_{i=1}^N \phi_{x_i} \otimes \psi_{y_i}$$

$$\hat{M}_{P_X \times P_Y} = \hat{M}_{P_X} \otimes \hat{M}_{P_Y}, \quad K_{ij} = k(x_i, x_j), \quad L_{ij} = l(y_i, y_j)$$

yielding the empirical HSIC = $\frac{1}{N^2} \text{Tr}(KHLH)$

where $H = I - \frac{1}{N} \mathbf{1}_N \mathbf{1}_N^T$

(Frobenius product of K, H, L, H)

$N \times N$ matrix of 1's

such that HKH is centered: has its column and row means subtracted

* In fact, this estimate is biased because it takes ~~over~~ products b/w repeated data points. To get an unbiased estimate we include only products b/w kernels of different data points.

A different approach to measuring dependence is to ~~minimize~~ find the pair of mappings $f: X \rightarrow \mathcal{F}, g: Y \rightarrow \mathcal{G}$ that maximizes the covariance b/w $f(x)$ and $g(y)$, $(x \sim P_X, y \sim P_Y)$ under a smoothness constraint on f and g . This is formalized as the Constrained Covariance (COCO):

$$\text{COCO}(P_{XY}, P_X, P_Y; \mathcal{F}, \mathcal{G}) = \sup_{\|f\|_{\mathcal{F}}=1, \|g\|_{\mathcal{G}}=1} \left(\mathbb{E}_{x,y \sim P_{XY}} f(x)g(y) - \mathbb{E}_{x \sim P_X} f(x) \mathbb{E}_{y \sim P_Y} g(y) \right)$$

In feature space,

$$= \sup_{\|f\|_{\mathcal{F}}=1, \|g\|_{\mathcal{G}}=1} \left(\langle M_{P_{XY}}, f \otimes g \rangle_{\mathcal{F} \otimes \mathcal{G}} - \langle \mu_{P_X}, f \rangle_{\mathcal{F}} \langle \mu_{P_Y}, g \rangle_{\mathcal{G}} \right)$$

$$= \sup_{\|f\|_{\mathcal{F}}=1, \|g\|_{\mathcal{G}}=1} \left(\langle M_{P_{XY}} - \mu_{P_X} \otimes \mu_{P_Y}, f \otimes g \rangle_{\mathcal{F} \otimes \mathcal{G}} \right)$$

$$= \sup_{\|f\|_{\mathcal{F}}=1, \|g\|_{\mathcal{G}}=1} \left(\langle C_{XY}, f \otimes g \rangle_{\mathcal{F} \otimes \mathcal{G}} \right)$$

C_{XY} (covariance)

for $f \in \mathcal{G}$
 $(L, (n \text{ obs}))_{\mathcal{F} \times \mathcal{G}}$
 $= (n, Lb)_{\mathcal{F}}$

$\sup_{\mathcal{F} \times \mathcal{G}} \langle f, C_{xy} g \rangle_{\mathcal{F} \times \mathcal{G}}$

\Rightarrow Empirical estimate of C_{xy} :

$$\hat{C}_{xy} = \frac{1}{N} \sum_{i=1}^N \phi_{xi} \otimes \psi_{yi} - (\hat{\mu}_x \otimes \hat{\mu}_y)$$

$$\hat{\mu}_x = \frac{1}{N} \sum_{i=1}^N \phi_{xi}, \quad \hat{\mu}_y \text{ analogous}$$

For $X = [\phi_{x1}, \dots, \phi_{xN}]$, $Y = [\psi_{y1}, \dots, \psi_{yN}]$

$$K = X^T X \text{ s.t. } K_{ij} = K(\phi_{xi}, \phi_{xj}), \quad L = Y^T Y \text{ s.t. } L_{ij} = L(\psi_{yi}, \psi_{yj})$$

$$\hat{C}_{xy} = \frac{1}{N} X H Y^T, \quad H = I - \frac{1}{N} \mathbb{1}_N$$

$$\hat{C}_{xy} = \left[\frac{1}{N} \sum_{i=1}^N \phi_{xi} \otimes \psi_{yi} \right] - \left(\frac{1}{N} \sum_{i=1}^N \phi_{xi} \right) \left(\frac{1}{N} \sum_{i=1}^N \psi_{yi} \right)$$

$$= \frac{1}{N} X Y^T - \left(\frac{1}{N} X \mathbb{1}_{N \times 1} \right) \left(\frac{1}{N} Y \mathbb{1}_{N \times 1} \right)$$

when $\mathbb{1}_{N \times 1}$ is an $N \times 1$ vector of 1's, such that $A \mathbb{1}_{N \times 1} = \text{sum of the columns of } A$

$$= \frac{1}{N} X Y^T - \frac{1}{N^2} (X \mathbb{1}_{N \times 1}) (Y \mathbb{1}_{N \times 1})^T$$

$$= \frac{1}{N} X Y^T - \frac{1}{N^2} X \mathbb{1}_{N \times N} Y^T$$

$$= \frac{1}{N} X \left(I - \frac{1}{N} \mathbb{1}_{N \times N} \right) Y^T$$

$$= \frac{1}{N} X H Y^T$$

We now solve the optimization problem above by assuming f, g are linear combinations of the data in feature space ~~where~~ $\{\tilde{\phi}_{xi}\}_{i=1}^N, \{\tilde{\psi}_{yi}\}_{i=1}^N$ respectively, ^{mean-subtracted} where $\tilde{\phi}_{xi} = \phi_{xi} - \frac{1}{N} \sum_{i=1}^N \phi_{xi}$

$$f = XH\alpha, \quad g = YH\beta$$

where again $H = I - \frac{1}{N} \mathbf{1}\mathbf{1}^T$ such that $XH = [\tilde{\phi}_{x1} \dots \tilde{\phi}_{xN}]$.

We then solve the following Lagrangian:

$$\mathcal{L}(f, g, \lambda_1, \lambda_2)$$

$$= f^T \tilde{C}_{xy} g - \frac{\lambda_1}{2} (\|f\|_2^2 - 1) - \frac{\lambda_2}{2} (\|g\|_2^2 - 1)$$

$$H = H^T \Rightarrow \alpha^T H X^T (X H Y^T) Y H \beta - \frac{\lambda_1}{2} (\alpha^T H X^T X H \alpha - 1) - \frac{\lambda_2}{2} (\beta^T H Y^T Y H \beta - 1)$$

$$= \frac{1}{N} \alpha^T H K H L H \beta - \frac{\lambda_1}{2} (\alpha^T H K H \alpha - 1) - \frac{\lambda_2}{2} (\beta^T H L H \beta - 1)$$

$$= \frac{1}{N} \alpha^T \tilde{K} \tilde{L} \beta - \frac{\lambda_1}{2} (\alpha^T \tilde{K} \alpha - 1) - \frac{\lambda_2}{2} (\beta^T \tilde{L} \beta - 1)$$

Taking the derivative w.r.t. α, β and setting to 0:

$$\frac{\partial}{\partial \alpha} \mathcal{L} = \frac{1}{N} \tilde{K} \tilde{L} \beta - \lambda_1 \tilde{K} \alpha = 0$$

$$\frac{\partial}{\partial \beta} \mathcal{L} = \frac{1}{N} \tilde{L} \tilde{K} \alpha - \lambda_2 \tilde{L} \beta = 0$$

If we multiply both sides by α^T for eqn 1 and by β^T for eqn 2,

$$\frac{1}{N} \alpha^T \tilde{K} \tilde{L} \beta = \lambda_1 \alpha^T \tilde{K} \alpha$$

$$\frac{1}{N} \beta^T \tilde{L} \tilde{K} \alpha = \lambda_2 \beta^T \tilde{L} \beta$$

Noting that $\alpha^T \tilde{K} \alpha = \|\tilde{K} \alpha\|_2^2 = 1$ and $\beta^T \tilde{L} \beta = \|\tilde{L} \beta\|_2^2 = 1$ and that $\alpha^T \tilde{K} \tilde{L} \beta = \text{Tr}(\alpha^T \tilde{K} \tilde{L} \beta) = \text{Tr}(\beta^T \tilde{L} \tilde{K} \alpha) = \text{Tr}(\beta^T \tilde{L} \tilde{K} \alpha) = \beta^T \tilde{L} \tilde{K} \alpha$, we have $\lambda_1 = \lambda_2$. Thus, we can solve for α, β by the following eigenvalue eqn:

$$\begin{bmatrix} 0 & \frac{1}{N} \tilde{K} \tilde{L} \\ \frac{1}{N} \tilde{L} \tilde{K} & 0 \end{bmatrix} \begin{bmatrix} \alpha \\ \beta \end{bmatrix} = \lambda \begin{bmatrix} \tilde{K} & 0 \\ 0 & \tilde{L} \end{bmatrix} \begin{bmatrix} \alpha \\ \beta \end{bmatrix}$$

where $\lambda = \lambda_1 = \lambda_2$

Effectively, the ~~smoothness~~ constraints $\|\tilde{K} \alpha\|_2 = 1$ and $\|\tilde{L} \beta\|_2 = 1$ enforce smoothness in f and g (recall that $\|\tilde{K} \alpha\|_2^2 = \sum_{\ell=-\infty}^{\infty} \frac{\hat{f}_\ell^2}{K_\ell}$), allowing only lower frequency ~~components~~ dependence to show up in the COCO. Higher frequency covariance can only be detected via COCO with large sample sizes.

Turns out $\text{HSIC} = \sum_{i=1}^N \gamma_i^2$, where γ_i is the i^{th} largest eigenvalue of $\begin{bmatrix} 0 & \frac{1}{N} \tilde{K} \\ \frac{1}{N} \tilde{L} & 0 \end{bmatrix}$ (from eigenvalue equation above), in the limit of infinite samples

(missing something here about how γ_i 's relate to f and g)

\tilde{C}_{xy} is defined as the matrix such that

$$\langle A, \tilde{C}_{xy} \rangle_{HS} = \mathbb{E}_{xy} \langle A, \phi(x) \otimes \psi(y) \rangle$$

μ_x is defined as the vector such that

$$\langle f, \mu_x \rangle = \mathbb{E}_x \langle f, x \rangle$$

○ We have $\|\hat{C}_{xy}\|^2 = \langle \hat{C}_{xy}, \hat{C}_{xy} \rangle$

$$= \frac{1}{N^2} \sum_{i,j} k(x_i, x_j) l(y_i, y_j)$$

$$= \frac{1}{N^2} \text{Tr}(KL)$$

Obviously, this is a biased estimate, since

○ we are including terms like $k(x_i, x_i), l(y_i, y_i)$

How biased is it?

Unbiased estimate: $\frac{1}{N(N-1)} \sum_{i=1}^N \sum_{j \neq i} k(x_i, x_j) l(y_i, y_j)$

Difference:

$$\frac{1}{N^2} \sum_{i,j} k(x_i, x_j) l(y_i, y_j) + \left(\frac{1}{N^2} - \frac{1}{N(N-1)} \right) \sum_{i=1}^N k(x_i, x_i) l(y_i, y_i)$$

$$= \frac{1}{N} \left[\frac{1}{N} \sum_{i=1}^N k(x_i, x_i) l(y_i, y_i) - \frac{1}{N(N-1)} \sum_{i=1}^N \sum_{j \neq i} k(x_i, x_j) l(y_i, y_j) \right]$$

Fidelity expectations:

$$= \left[\mathbb{E}_{xy} K(x, x) \ell(y, y) - \mathbb{E}_{xy} \mathbb{E}_{x'y'} K(x, x') \ell(y, y') \right] \frac{1}{n}$$

\Rightarrow i.e. the ^{expected} difference/drops with $\frac{1}{n}$ (BIAS)

(*but, be careful about kernel such that $K(x, x), \ell(y, y)$ not too large such that they dominate)

Statistical Testing w/ HSIC

The (biased) empirical estimate of HSIC

\Rightarrow a χ^2 -statistic when $P_{xy} \neq P_x P_y$

Under the null, however, ($P_{xy} = P_x P_y$), ~~it's~~

it is a ~~degenerate~~ degenerate χ^2 -statistic

But we can find

^{1st and 2nd} the moments, and approximate

the infinite sum of χ^2 's w/ a Gamma distribution with matched first two moments.

Or, we can use a permutation test.

Hilbert-Schmidt Operators

~~Let~~ $L, M: G \rightarrow F$, i.e. $L, M \in F \times G$

Let $\{f_i\}_{i \in I}$ be basis of F
 $\{g_j\}_{j \in J}$ " " " " G

Dot Product:

$$\begin{aligned} \langle L, M \rangle_{HS} &= \sum_{j \in J} \langle Lg_j, Mg_j \rangle_F \\ &= \sum_{i \in I} \sum_{j \in J} \langle Lg_j, f_i \rangle_F \langle Mg_j, f_i \rangle_F \end{aligned}$$

Norm:

$$\|L\|_{HS}^2 = \sum_{j \in J} \|Lg_j\|_F^2 = \sum_{i \in I} \sum_{j \in J} |\langle Lg_j, f_i \rangle_F|^2$$

Rank 1 operators: $(b \otimes a)g \rightarrow \langle g, a \rangle_F b$

for $a \in F$, $b \in G$ forming rank 1 operator $a \otimes b: G \rightarrow F$

$$\begin{aligned} \text{Property \# 1: } \|a \otimes b\|_{HS}^2 &= \sum_{g \in G} \|(a \otimes b)g\|_F^2 \\ &= \sum_{g \in G} \|\langle g, b \rangle_G a\|_F^2 \\ &= \|a\|_F^2 \sum_{g \in G} |\langle g, b \rangle_G|^2 \\ &= \|a\|_F^2 \|b\|_G^2 \end{aligned}$$

Property # 2: $\langle L, a \otimes b \rangle_{HS} = \langle a, Lb \rangle_{\mathcal{F}}$

$$\text{Pf. } \langle L, a \otimes b \rangle_{HS} = \sum_{j \in S} \langle Lg_j, (a \otimes b)g_j \rangle_{\mathcal{F}}$$

$$= \sum_{j \in S} \langle Lg_j, a \langle b, g_j \rangle_{\mathcal{G}} \rangle_{\mathcal{F}}$$

$$= \sum_{j \in S} \langle Lg_j, a \rangle_{\mathcal{F}} \langle b, g_j \rangle_{\mathcal{G}}$$

$$= \left\langle \sum_{j \in S} L \langle b, g_j \rangle_{\mathcal{G}} g_j, a \right\rangle_{\mathcal{F}}$$

$$= \left\langle L \sum_{j \in S} \langle b, g_j \rangle_{\mathcal{G}} g_j, a \right\rangle_{\mathcal{F}}$$

$$= \underline{\underline{\langle Lb, a \rangle_{\mathcal{F}}}}$$

Property # 3: $\langle u \otimes v, a \otimes b \rangle_{HS} = \langle u, a \rangle_{\mathcal{F}} \langle v, b \rangle_{\mathcal{G}}$

$$\text{Pf. From above, } \langle u \otimes v, a \otimes b \rangle_{HS} = \langle a, (u \otimes v)b \rangle_{\mathcal{F}}$$

$$= \langle a, \langle v, b \rangle_{\mathcal{G}} u \rangle_{\mathcal{F}}$$

$$= \underline{\underline{\langle a, u \rangle_{\mathcal{F}} \langle v, b \rangle_{\mathcal{G}}}}$$