

---

# Supplementary material for "A Kernel Two Sample Test for Unbiased Decisions"

---

Anonymous Author(s)

Affiliation

Address

email

1 **Outline.** This document provides supplementary material accompanying the main body of the  
2 paper "A Kernel Two-Sample Test for Unbiased Decisions". It includes in section A proofs of all  
3 propositions and theorems; in section B a more detailed description of the example provided in the  
4 introduction; in section C a detailed description of other tests and our implementations; and finally, in  
5 section D a discussion of computational complexity and possible methods to speed up computations.

## 6 A. Proofs

7 In this section we prove the propositions and theorems described in the main body of this paper.

### 8 A.1. Proof of proposition 1

9 Assume that kernel  $k : \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}$  is characteristic and that for all  $y$ ,  $w(x) > 0$  is bounded above  
10 by  $W$ . A kernel is called characteristic, if the maximum mean discrepancy between probability  
11 measures  $P_{Y^0}$  and  $P_{Y^1}$  induced by  $k$  is such that,  $\text{MMD}(P_{Y^0}, P_{Y^1}) = 0$  if and only if  $P_{Y^0} = P_{Y^1}$ .  
12 [2] showed that Gaussian kernels are characteristic.

13 To prove the proposition we exploit the assumption  $Y^0, Y^1 \perp\!\!\!\perp T | X$  and recover expectations with  
14 respect to the underlying random variables of interest ( $Y^0, Y^1$ ). Assuming access to the propensity  
15 score,  $e(x) = p(T = 1 | X = x) = \mathbb{E}(I(T = 1) | X = x)$ , and for any measurable function of our  
16 observed values  $Y$ , such as the kernel function  $k$ , we have that,

$$\begin{aligned} \mathbb{E}_{y, y^* \sim \mathcal{Y} | T=1} \left( \frac{k(y, y^*)}{e(X)e(X^*)} \right) &= \mathbb{E}_{y, y^*} \left( \frac{I(T=1)I(T^*=1)k(y, y^*)}{e(X)e(X^*)} \right) \\ &= \mathbb{E}_{y, y^*} \left( \frac{I(T=1)I(T^*=1)k(y^1, y^{1*})}{e(X)e(X^*)} \right) \\ &= \mathbb{E}_{x, x^*} \left( \mathbb{E}_{y, y^*} \left( \frac{I(T=1)I(T^*=1)k(y^1, y^{1*})}{e(x)e(x^*)} \mid y^1, y^{1*}, x, x^* \right) \right) \\ &= \mathbb{E}_{y^1, y^{1*}, x, x^*} \left( \frac{k(y^1, y^{1*})}{e(x)e(x^*)} \mathbb{E}_{T, T^*} (I(T=1)I(T^*=1) | y^1, y^{1*}, x, x^*) \right) \\ &= \mathbb{E}_{y^1, y^{1*}} (k(y^1, y^{1*})) \end{aligned}$$

17 where recall that we use the notation  $y^1$  for a realization of the random variable  $Y^1$ .  $I$  is the  
18 indicator function. This derivation shows that by taking weighted expectations with respect to  
19 the observed distribution  $Y | T = 1$  we can access expectations with respect to our distribution  
20 of interest  $Y^1$ . Similar derivations follow for data observed under  $Y | T = 0$  using the fact that

21  $\mathbb{E}_{Y | T=0} \left( \frac{f(Y)}{1-e(X)} \right) = \mathbb{E}_{Y^0} (f(Y^0))$ , for  $f$  any measurable function.

22 Now notice that the  $\text{MMD}(Y^0, Y^1)$  between  $Y^0$  and  $Y^1$  is defined in terms of expectations with  
 23 respect to the random variables  $Y^0$  and  $Y^1$ ,

$$\text{MMD}(Y^0, Y^1) := \mathbb{E}_{y^0, y^{0,*}} k(y^0, y^{0,*}) + \mathbb{E}_{y^1, y^{1,*}} k(y^1, y^{1,*}) - 2\mathbb{E}_{y^1, y^0} k(y^0, y^1)$$

24 Thus with the above derivation we get that each term in the definition of  $\text{WMMD}(Y|T=0, Y|T=1)$   
 25 is equal to each term in the definition of the MMD, which proves the proposition.  $\square$

## 26 A.2. Proof of Theorem 1

27 **Regularity conditions.** The following notation is used in the statement on the regularity conditions  
 28 of Theorem 1. Let  $B_n = (b_{imn})$  and  $W_n = (W_{ijn})$ , for  $i, j = 1, \dots, n; n, m : 1, 2, \dots$ . Here  $W_n$  is  
 29 a matrix of weights in  $\mathbb{R}^{n \times n}$  and  $B_n$  is an orthogonal matrix in  $\mathbb{R}^{m \times n}$  such that  $B_n^T W_n B_n = \Lambda_n$ ,  
 30 where  $\Lambda_n$  is a diagonal matrix with  $\lambda_{mn}$  as the  $m^{\text{th}}$  diagonal element. Assume  $\lim_{n \rightarrow \infty} \lambda_{mn} = \lambda_m$   
 31 and let  $\delta_{km}$  be the dirac delta function with  $\delta_{km} = 1$  if  $k = m$  and zero otherwise. Assume that the  
 32 following regularity conditions hold,

- 33 1.  $\max_{1 \leq i \leq n} |b_{imn}| \rightarrow 0$  as  $n \rightarrow \infty$  for each  $m$ .
- 34 2.  $\sum_{i=1}^n b_{imn} b_{ikn} \rightarrow \delta_{mk}$  as  $n \rightarrow \infty$  for all  $m, k$ .
- 35 3.  $\sum_{i=1}^n \sum_{j=1}^n w_{ijn}^2 \rightarrow \sum_{m=1}^{\infty} \lambda_m^2 < \infty$ .
- 36 4.  $\sum_{i=1}^n \sum_{j=1}^n w_{ijn} b_{ikn} b_{jkn} \rightarrow \lambda_k$  as  $n \rightarrow \infty$ , for all  $m$ .

37 These conditions are sufficient by [1] for a square matrix of data-dependent weights  $W = (w_{ij})$  to  
 38 be approximately diagonalizable, such that it admits an eigen-decomposition  $B^T W B = \Lambda$ .

39 **Proof.** Recall the definition of the empirical estimate of the  $\text{WMMD}^2$ ,

$$\widehat{\text{WMMD}}^2 := \frac{1}{n(n-1)} \sum_{i \neq j: t_i = t_j = 1} w_i w_j k(y_i, y_j) + \frac{1}{m(m-1)} \sum_{i \neq j: t_i = t_j = 0} k(y_i, y_j) - \quad (1)$$

$$\frac{2}{nm} \sum_{i, j: t_i = 1, t_j = 0} w(x_i) k(y_i, y_j) \quad (2)$$

40 where the  $(y_i, t_i, x_i)$  are realization of the random variables  $(Y, T, X)$ , and have assumed that  $n$   
 41 observations are made with  $T = 1$  and  $m$  with  $T = 0$ .  $w(x_i) = \text{Pr}(T_i = 1 | X_i = x_i) / \text{Pr}(T_i =$   
 42  $0 | X_i = x_i)$  is the density ratio giving the likelihood of an example  $i$  being observed under one  
 43 population with respect to the other. We assume this ratio to be known (for now) and provide  
 44 approximation bounds for our proposed approximation in Theorem 2 and 3. Our proof is presented in  
 45 three parts, each one deriving the asymptotic behaviour of each one of the three terms in (1).

46 Note first that we may write the square integrable (centered) kernel  $k$  as a weighted sum of product  
 47 of eigen-functions of the Hilbert-Schmidt operator defined by  $k$  [3],

$$k(y_i, y_j) = \sum_{k=1}^{\infty} \alpha_k \psi_k(Y_i) \psi_k(Y_j) \quad (3)$$

48 Consider now the first term in (1), it follows that,

$$\sum_{i=1}^n \sum_{j=1, j \neq i}^n w(x_i) w(x_j) k(y_i, y_j) = \sum_{i=1}^n \sum_{j=1, j \neq i}^n w_{ij} \sum_{k=1}^{\infty} \alpha_k \psi_k(Y_i) \psi_k(Y_j) \quad (4)$$

$$= \sum_{k=1}^{\infty} \alpha_k \sum_{i=1}^n \sum_{j=1, j \neq i}^n w_{ij} \psi_k(Y_i) \psi_k(Y_j) \quad (5)$$

49 where we have dropped the  $t_i$ 's in the summation indices and have written  $w_{ij} = w(x_i)w(x_j)$  for  
 50 brevity. Using the degeneracy of  $k$  (in the sense that  $\text{Var}[\mathbb{E}[k(y, y')]] = 0$ ), the eigen-functions  
 51  $\psi_k(Y_i)$ ,  $i = 1, \dots, n$  are zero mean independent random variables by the independence of the  $Y_i$ .  
 52 Using the above and the regularity conditions, Theorem 1 in [14] yields,

$$\frac{1}{n} \sum_{i=1}^n \sum_{j=1, j \neq i}^n w_{ij} \psi_k(Y_i) \psi_k(Y_j) \xrightarrow{d} \sum_{m=1}^{\infty} \lambda_m (Z_{km}^2 - 1) \quad (6)$$

53 where  $Z_{km} \sim \mathcal{N}(0, 1)$  are *i.i.d.*.

54 The limiting distribution of the un-weighted term in (1) is that of a well-studied U-Statistic whose  
55 derivation can be found in Section 5.5.2 of [10].

$$\frac{1}{m} \sum_{i=1:t_i=0}^m \sum_{j=1:j \neq i:t_j=0}^m k(Y_i, Y_j) \xrightarrow{d} \sum_{k=1}^{\infty} \alpha_k (V_k^2 - 1) \quad (7)$$

56 The limiting distribution of the cross term in (1) follows from a modification of the derivation of  
57 Theorem 1 in [1] and is given by,

$$\frac{1}{\sqrt{nm}} \sum_{i=1:t_i=1}^n \sum_{j=1:t_j=0}^n w'_{ij} \psi_k(Y_i) \psi_k(Y_j) \xrightarrow{d} \sum_{m=1}^{\infty} \lambda'_m Z_{km} V_{km} \quad (8)$$

58 where the eigenvalues  $(\lambda'_m)$  correspond to those of the eigen-decomposition of the weight matrix  
59  $W'$  with  $W'_{ij} = w(x_i)$  and where  $V_{km} \sim \mathcal{N}(0, 1)$  independently of  $Z_{km} \sim \mathcal{N}(0, 1)$ . We prove (8)  
60 below.

61 We now combine these results. Define  $t = m + n$ , and assume  $\lim_{m,n \rightarrow \infty} m/t = \rho_y$  and  
62  $\lim_{m,n \rightarrow \infty} n/t = \rho_x := (1 - \rho_y)$  for fixed  $0 < \rho_x < 1$ . Then,

$$t \widehat{\text{WMMD}}^2 \xrightarrow{d} \rho_x^{-1} \sum_{k=1}^{\infty} \alpha_k \sum_{m=1}^{\infty} \lambda_m (Z_{km}^2 - 1) + \rho_y^{-1} \sum_{k=1}^{\infty} \alpha_k (V_k^2 - 1) - \frac{2}{\sqrt{\rho_x \rho_y}} \sum_{k=1}^{\infty} \alpha_k \sum_{m=1}^{\infty} \lambda'_m Z_{km} V_{km} \quad (9)$$

63 In the case that both samples have equal size with total sample size  $n$ , we have that under  $\mathcal{H}_0$ ,

$$n \widehat{\text{WMMD}}^2 \xrightarrow{d} \sum_{k=1}^{\infty} \alpha_k \sum_{m=1}^{\infty} \lambda_m (Z_{km}^2 - 1) + \sum_{k=1}^{\infty} \alpha_k (V_k^2 - 1) - 2 \sum_{k=1}^{\infty} \alpha_k \sum_{m=1}^{\infty} \lambda'_m Z_{km} V_{km} \quad (10)$$

64 **The case of  $P \neq Q$ , under  $\mathcal{H}_1$ .** The centered kernel  $k$  is non-degenerate since its expectation  
65 under assumption  $\mathcal{H}_1$  is different from 0. The limiting distribution of WMMD can be derived by  
66 considering each term in the sum separately. For the first and third terms,

$$(\star) := \frac{1}{n(n-1)} \sum_{i \neq j: t_i = t_j = 1} w(x_i) w(x_j) k(y_i, y_j), \quad (\star\star) := \frac{2}{mn} \sum_{i, j: t_i = 1, t_j = 0} w(x_i) k(y_i, y_j) \quad (11)$$

67 we get immediately by Theorem 2.1 from p. 4, [11] that their limiting distributions are normal with  
68 mean  $\mathbb{E}(\star)$  and variance  $\text{Var}(\star)$ , and mean  $\mathbb{E}(\star\star)$  and variance  $\text{Var}(\star\star)$ , respectively. The middle term  
69  $\frac{1}{m(m-1)} \sum_{i \neq j: t_i = t_j = 0} k(y_i, y_j)$  is an un-weighted U-statistic whose limiting distribution is given by  
70 the results in section 5.5 [10]. As above, define  $t = m + n$ , and assume  $\lim_{m,n \rightarrow \infty} m/t = \rho_y$  and  
71  $\lim_{m,n \rightarrow \infty} n/t = \rho_x := (1 - \rho_y)$  for fixed  $0 < \rho_x < 1$ . Collecting these results, we get under  $\mathcal{H}_1$ ,

$$t^{1/2} \left( \widehat{\text{WMMD}}^2 - \text{WMMD}^2 \right) \xrightarrow{d} \mathcal{N}(0, \sigma_{\mathcal{H}_1}^2) \quad (12)$$

72 where we write  $z = ((y_1, t = 1, x_1), (y_0, t = 0, x_0))$  for the joint sample under the two populations,  
73 and  $h(z, z^*) := w(x_1) w(x_1^*) k(y_1, y_1^*) + \mathbb{E} k(y_0, y_0^*) - 2w(x_1) k(y_1, y_0^*)$ .  $\sigma_{\mathcal{H}_1}^2 := \text{Var}_z(\mathbb{E}_{z^*} h(z, z^*))$   
74 [10, 3].

75 **Proof of equation (8).** The proof is a modification of the result of the convergence of degenerate U  
76 statistics on p. 761 in [3] and of the derivation of Theorem 1 in [1].

77 Consider,

$$T_k := \frac{1}{\sqrt{nm}} \sum_{i=1:t_i=1}^n \sum_{j=1:t_j=0}^m w'_{ij} \psi_k(Y_i) \psi_k(Y_j) \quad (13)$$

78 and define for each  $k$ ,

$$w_{ij}^* := \sum_{s=1}^S \lambda_s b_{isk} b_{jsk}, \quad T_k^* := \frac{1}{\sqrt{nm}} \sum_{i=1:t_i=1}^n \sum_{j=1:t_j=0}^n w'_{ij} \psi_k(Y_i) \psi_k(Y_j) \quad (14)$$

79 We will start by showing that  $\sum_{i=1}^n \sum_{j=1}^m (w_{ij} - w_{ij}^*)^2 \rightarrow 0$  as  $n, m \rightarrow \infty$ . Note that this implies  
80 that  $\text{Var}(T_k^* - T_k) \rightarrow 0$  and thus that the distributions of  $T_k^*$  and  $T_k$  coincide in the limit. We will  
81 proceed by showing first the convergence of the sum of squares and then we derive the distribution of  
82  $T_k^*$ . Using the definitions above, write,

$$\begin{aligned}
\sum_{i=1}^n \sum_{j=1}^m (w_{ij} - w_{ij}^*)^2 &= \sum_{i=1}^n \sum_{j=1}^m w_{ij}^2 - 2 \sum_{s=1}^S \lambda_s \sum_{i=1}^n \sum_{j=1}^m w_{ij} b_{isk} b_{jks} + \\
&\quad \sum_{s=1}^S \sum_{t=1}^S \lambda_s \lambda_t \left( \sum_{i=1}^n b_{isk} b_{itk} \right) \left( \sum_{j=1}^m b_{jks} b_{jtk} \right) \\
&= \sum_{i=1}^n \sum_{j=1}^m w_{ij}^2 - \sum_{s=1}^S \lambda_s^2 - 2 \sum_{s=1}^S \lambda_s \left( \sum_{i=1}^n \sum_{j=1}^m w_{ij} b_{isk} b_{jks} - \lambda_s \right) \\
&\quad + \sum_{s=1}^S \sum_{t=1}^S \lambda_t \lambda_s \left( \sum_{i=1}^n b_{isk} b_{itk} - \delta_{st} \right) \left( \sum_{j=1}^m b_{jks} b_{jtk} - \delta_{st} \right) + \\
&\quad 2 \sum_{s=1}^S \lambda_s^2 \left( \sum_{i=1}^n \sum_{j=1}^m b_{iks}^2 - 1 \right) \tag{15}
\end{aligned}$$

83 where we have removed the group allocation indices  $t$  for clarity. Note here that the first and second  
84 term cancel each other by Assumption 1 of the regularity conditions, the third term is  $\mathcal{O}(1)$  by  
85 Assumption 4 and the fourth and fifth terms are also  $\mathcal{O}(1)$  by Assumption 2 and the properties of the  
86 dirac delta function.

87 Consider now  $T_k^*$  and rewrite it as,

$$T_k^* = \sum_{s=1}^S \lambda_s \left( \frac{1}{\sqrt{n}} \sum_{i=1:t_i=1}^n b_{isk} \psi_k(Y_i) \right) \left( \frac{1}{\sqrt{m}} \sum_{j=1:t_j=0}^m b_{jks} \psi_k(Y_j) \right) \tag{16}$$

88 Define the length  $K$  vectors  $\Psi_n$  and  $\Psi'_m$  having  $k_{th}$  entries,

$$\Psi_{kn} = \left( \frac{1}{\sqrt{n}} \sum_{i=1:t_i=1}^n b_{isk} \psi_k(Y_i) \right), \quad \Psi'_{km} = \left( \frac{1}{\sqrt{m}} \sum_{j=1:t_j=0}^m b_{jks} \psi_k(Y_j) \right) \tag{17}$$

89 respectively. These have mean and covariance,

$$\mathbb{E}(\Psi_{kn}) = 0, \quad \text{Cov}(\Psi_{kn}, \Psi_{k'n}) = \begin{cases} \frac{1}{m} \sum_{i=1}^n b_{isk}^2 = 1, & \text{if } k = k' \\ 0, & \text{otherwise} \end{cases} \tag{18}$$

90 Moreover, the vectors  $\Psi_n$  and  $\Psi'_m$  are independent. The results (8) then holds by the Lindberg-Levy  
91 Central Limit Theorem [10], Theorem 1.9.1A.  $\square$

### 92 A.3. Proof of Theorem 2

93 We assume that for increasing sample size, as  $n, m \rightarrow \infty$ , we can approximate arbitrarily well the  
94 density ratio  $w(x)$ , for all  $x$  in our training data. This is justified by the following Lemma,

95 **Lemma 1** (Lemma 1.4 [5]) *Let  $w(x_i) \in [0, B]$  be the optimal weight in the population sense,*  
96  *$\Pr(T_i = 1|x_i) = w(x_i)\Pr(T_i = 0|x_i)$ . Assume we draw  $n$  samples from  $X|T = 1$  and  $m$  samples*  
97 *from  $X|T = 1$  independently and that  $\|\phi(x)\| \leq R$ . Then, with probability at least  $1 - \delta$ ,*

$$\left\| \frac{1}{n} \sum_{i=1:t_i=1}^n w(x_i) \phi(x_i) - \frac{1}{m} \sum_{j=1:t_j=0}^m \phi(x_j) \right\| \leq \left( 1 + \sqrt{2 \log \frac{2}{\delta}} \right) R \sqrt{\frac{B^2}{n} + \frac{1}{m}} \tag{19}$$

98 Note that because the optimization problem is convex the choice of  $\hat{w}(x) := \Pr(T_i = 1|x)/\Pr(T_i =$   
99  $0|x)$  uniquely minimizes the objective function with value 0, see Lemma 1.3, [5]. Thus by the

100 argument above, we may assume that for increasing sample size, as  $n, m \rightarrow \infty$ ,  $\hat{w}(x) \rightarrow w(x)$ , for  
 101 all  $x$  in the common support of the distributions  $Pr(T_i = 1|x)$  and  $Pr(T_i = 0|x)$ .

102 Consider the first terms of  $\widehat{\text{WMMD}}^2(\hat{w})$  and  $\widehat{\text{WMMD}}^2(w)$ , that denote the empirical WMMD<sup>2</sup>  
 103 with estimated and true weights  $w$  respectively,

$$\hat{K}_{n,m} := \sum_{i=1:t_i=1}^n \sum_{j=1, j \neq i:t_j=1}^m \hat{w}_{ij} k(y_i, y_j), \quad \text{and} \quad K_{n,m} := \sum_{i=1:t_i=1}^n \sum_{j=1, j \neq i:t_j=1}^m w_{ij} k(y_i, y_j) \quad (20)$$

104 It holds that  $\sum_{i=1}^n \sum_{j=1, j \neq i}^m (\hat{w}_{ij} - w_{ij})^2 \rightarrow 0$  as  $n, m \rightarrow \infty$  by the arguments at the end of section  
 105 A.3. This implies that  $\text{Var}(\hat{K}_{n,m} - K_{n,m}) \rightarrow 0$  and  $E(|\hat{K}_{n,m} - K_{n,m}|^2) \rightarrow 0$  which means that  
 106  $\hat{K}_{n,m} - K_{n,m}$  converges to 0 in  $L_2$ , and hence in distribution. The distributions of  $\hat{K}_{n,m}$  and  $K_{n,m}$   
 107 coincide in the limit.

108 The same derivations apply for the other two terms in the definition of  $\widehat{\text{WMMD}}^2$ . Therefore we con-  
 109 clude that  $\widehat{\text{WMMD}}^2$  with estimated weights has the same asymptotic null and alternative distribution  
 110 as  $\widehat{\text{WMMD}}^2$  with known weights. In particular, asymptotically, its false positive rate is  $\alpha$  and its  
 111 power converges to 1.

### 112 A.5. Proof of Theorem 3

113 We prove Theorem 3 by first stating and proving several Lemmas which bound the different terms of  
 114 the inequality of interest.

115 **Lemma 2** *In addition to the conditions of Lemma 1, assume there exists some  $\hat{w}_i$ , the empirical*  
 116 *counterparts of the population weights estimated by matching kernel mean embeddings, such that,*

$$\left\| \frac{1}{n} \sum_{i=1:t_i=1}^n \hat{w}_i \phi(x_i) - \frac{1}{m} \sum_{j=1:t_j=0}^m \phi(x_j) \right\| \leq \epsilon \quad (21)$$

117 *Then,*

$$\left\| \frac{1}{n} \sum_{i=1}^n w_i \phi(x_i) - \frac{1}{n} \sum_{i=1}^n w(x_i) \phi(x_i) \right\| \leq \epsilon + \left( 1 + \sqrt{2 \log \frac{2}{\delta}} \right) R \sqrt{\frac{B^2}{n} + \frac{1}{m}} \quad (22)$$

118 *Proof.* Note that by using Lemma 1 and the triangle inequality we immediately get,

$$\begin{aligned} \left\| \frac{1}{n} \sum_{i=1:t_i=1}^n \hat{w}_i \phi(x_i) - \frac{1}{n} \sum_{i=1:t_i=1}^n w_i \phi(x_i) \right\| &\leq \left\| \frac{1}{n} \sum_{i=1:t_i=1}^n \hat{w}_i \phi(x_i) - \frac{1}{m} \sum_{j=1:t_j=0}^m \phi(x_j) \right\| \\ &+ \left\| \frac{1}{n} \sum_{i=1:t_i=1}^n w_i \phi(x_i) - \frac{1}{m} \sum_{j=1:t_j=0}^m \phi(x_j) \right\| \\ &\leq \epsilon + \left( 1 + \sqrt{2 \log \frac{2}{\delta}} \right) R \sqrt{\frac{B^2}{n} + \frac{1}{m}} \end{aligned} \quad (23)$$

119 □

120 **Lemma 3** *Let  $\widehat{\text{WMMD}}(w)$  be the weighted estimator of the MMD given i.i.d. distorted samples as*  
 121 *defined in (1) with known (population) weights  $w$ , and similarly define  $\widehat{\text{WMMD}}(\hat{w})$  with weights  $\hat{w}$*   
 122 *estimated by matching the empirical kernel mean embeddings of the distorted samples. Then, given*  
 123 *the conditions of Lemmas 1 and 2,*

$$\left| \widehat{\text{WMMD}}^2(\hat{w}) - \widehat{\text{WMMD}}^2(w) \right| \leq 2R(B+1) \left( \epsilon + \left( 1 + \sqrt{2 \log \frac{2}{\delta}} \right) R \sqrt{\frac{B^2}{n} + \frac{1}{m}} \right) \quad (24)$$

124 *Proof.* Consider expanding the estimators,

$$\begin{aligned} \left| \widehat{\text{WMMD}}^2(\hat{w}) - \widehat{\text{WMMD}}^2(w) \right| &= \frac{1}{n(n-1)} \sum_{i \neq j} \hat{w}_i \hat{w}_j k(y_i, y_j) - \frac{1}{n(n-1)} \sum_{i \neq j} w_i w_j k(y_i, y_j) \\ &\quad - \left( \frac{2}{nm} \sum_{i,j} \hat{w}_i k(y_i, y_j) - \frac{2}{nm} \sum_{i,j} w(x_i) k(y_i, y_j) \right) \end{aligned} \quad (25)$$

125 Note that the U-statistic in  $y$  cancel since these do not involve the weights.

126 **First and second terms.** We can bound the first and second terms as follows,

$$\frac{1}{n(n-1)} \sum_{i \neq j} \hat{w}_i \hat{w}_j k(y_i, y_j) - \frac{1}{n(n-1)} \sum_{i \neq j} w_i w_j k(y_i, y_j) \quad (26)$$

$$= \frac{1}{n(n-1)} \sum_{i \neq j} \hat{w}_i \hat{w}_j \langle \psi(y_i), \psi(y_j) \rangle - \frac{1}{n(n-1)} \sum_{i \neq j} w_i w_j \langle \psi(y_i), \psi(y_j) \rangle \quad (27)$$

$$\begin{aligned} &= \left| \left\langle \frac{1}{n} \sum_{i=1}^n \hat{w}_i \psi(y_i) - \frac{1}{n} \sum_{i=1}^n w(x_i) \psi(y_i), \frac{1}{n-1} \sum_{j=1, j \neq i}^m \hat{w}_j \psi(y_j) \right\rangle \right. \\ &\quad \left. + \left\langle \frac{1}{n} \sum_{i=1}^n \hat{w}_i \psi(y_i) - \frac{1}{n} \sum_{i=1}^n w(x_i) \psi(y_i), \frac{1}{n-1} \sum_{j=1, j \neq i}^m w(x_j) \psi(y_j) \right\rangle \right| \end{aligned} \quad (28)$$

$$\begin{aligned} &\leq \left| \left\langle \frac{1}{n} \sum_{i=1}^n \hat{w}_i \psi(y_i) - \frac{1}{n} \sum_{i=1}^n w(x_i) \psi(y_i), \frac{1}{n-1} \sum_{j=1, j \neq i}^m \hat{w}_j \psi(y_j) \right\rangle \right| \\ &\quad + \left| \left\langle \frac{1}{n} \sum_{i=1}^n \hat{w}_i \psi(y_i) - \frac{1}{n} \sum_{i=1}^n w(x_i) \psi(y_i), \frac{1}{n-1} \sum_{j=1, j \neq i}^m w(x_j) \psi(y_j) \right\rangle \right| \end{aligned} \quad (29)$$

$$\leq 2BR \left( \epsilon + \left( 1 + \sqrt{2 \log \frac{2}{\delta}} \right) R \sqrt{\frac{B^2}{n} + \frac{1}{m}} \right) \quad (30)$$

127 where  $\psi(\cdot) := k(y, \cdot)$ . Note that we have omitted the group allocation indices, these should be  
 128 clear however from the  $i$  and  $j$  indices. The second equality follows by adding and subtracting  
 129  $\frac{1}{n(n-1)} \sum_{i=1}^n \sum_{j=1, j \neq i}^m w(x_i) \hat{w}_j \langle \psi(y_i), \psi(y_j) \rangle$  which factorizes into the given expression. The  
 130 second to last inequality follows from the triangle inequality and the last inequality follows from the  
 131 properties of norms and the results derived in Lemmas 1 and 2.

132 **Third and fourth terms.** The third and fourth terms (in brackets) are derived similarly and satisfy  
 133 the following bounds,

$$\frac{2}{nm} \sum_{i,j} \hat{w}_i k(y_i, y_j) - \frac{2}{nm} \sum_{i,j} w(x_i) k(y_i, y_j) \quad (31)$$

$$= \frac{2}{nm} \sum_{i,j} \hat{w}_i \langle \psi(y_i), \psi(y_j) \rangle - \frac{2}{nm} \sum_{i,j} w(x_i) \langle \psi(y_i), \psi(y_j) \rangle \quad (32)$$

$$= \left| \frac{1}{n} \sum_{i=1}^n \hat{w}_i \left\langle \psi(y_i), \frac{2}{m} \sum_{j=1}^m \psi(y_j) \right\rangle - \frac{1}{n} \sum_{i=1}^n w(x_i) \left\langle \psi(y_i), \frac{2}{m} \sum_{j=1}^m \psi(y_j) \right\rangle \right| \quad (33)$$

$$= \left| \left\langle \frac{1}{n} \sum_{i=1}^n \hat{w}_i \psi(y_i) - \frac{1}{n} \sum_{i=1}^n w(x_i) \psi(y_i), \frac{2}{m} \sum_{j=1}^m \psi(y_j) \right\rangle \right| \quad (34)$$

$$\leq 2R \left( \epsilon + \left( 1 + \sqrt{2 \log \frac{2}{\delta}} \right) R \sqrt{\frac{B^2}{n} + \frac{1}{m}} \right) \quad (35)$$

134 where the last inequality follows from the properties of norms and the results derived in Lemmas 1  
135 and 2.

136 Finally, collecting the two bounds the lemma follows.  $\square$

137 **Lemma 4** Let  $\widehat{\text{WMMD}}(w)$  be the weighted estimator of the MMD given *i.i.d.* distorted samples  
138 as defined in (1) with known (population) weights  $w$ , and maximum kernel value  $R$ . Assume that  
139  $1 \leq w \leq B$  for all  $x \in \mathcal{X}$ . Then, with probability at least  $1 - \delta$ ,

$$\left| \widehat{\text{WMMD}}^2(w) - \text{MMD}^2 \right| \leq R(B+1)^2 \sqrt{\frac{1}{2m_2} \log \frac{1}{\delta}} \quad (36)$$

140 where  $m_2 := \lfloor m/2 \rfloor$ .

141 *Proof.* Assuming the kernel  $k(\cdot, \cdot)$  is bounded between 0 and  $R$  and the weights  $w$  bounded between 0  
142 and  $B$ , we can infer function bounds such that  $-2BR \leq wk(y_i, x_j) \leq R(B^2 + 1)$ . By Theorem 10  
143 in [3] which results from an application of the large deviation bound on U statistics due to Hoeffding  
144 we have that,

$$p \left( \left| \widehat{\text{WMMD}}^2(w) - \text{MMD}^2 \right| > e \right) \leq \exp \left\{ \frac{-2e^2 m_2}{R^2(B+1)^4} \right\} \quad (37)$$

145 Define  $\delta = \exp \left\{ \frac{-2e^2 m_2}{R^2(B+1)^4} \right\}$ . Thus, with probability  $1 - \delta$ ,

$$\left| \widehat{\text{WMMD}}^2(w) - \text{MMD}^2 \right| \leq R(B+1)^2 \sqrt{\frac{1}{2m_2} \log \frac{1}{\delta}} \quad (38)$$

146 where  $m_2 := \lfloor m/2 \rfloor$ .  $\square$

147 We are ready to prove Theorem 3. This will be a straightforward combination of the lemmas given  
148 above.

149 **Proof of Theorem 3.** Let  $\widehat{\text{WMMD}}(\hat{w})$  be the weighted estimator of the MMD given *i.i.d.* distorted  
150 samples as defined in (1) with estimated weights  $\hat{w}$ . Assume conditions on Lemmas 1,2,3 and 4  
151 above hold and that there exists an  $\epsilon > 0$  such that,

$$\left\| \frac{1}{n} \sum_{i=1}^n \hat{w}_i \phi(x_i) - \frac{1}{m} \sum_{i=1}^m \phi(x_i) \right\| \leq \epsilon \quad (39)$$

152 We may decompose the absolute difference between our weighted approximation using distorted  
153 samples and the population MMD as follows,

$$\begin{aligned} & \left| \widehat{\text{WMMD}}^2(\hat{w}) - \text{MMD}^2 \right| \\ & \leq \left| \widehat{\text{WMMD}}^2(\hat{w}) - \widehat{\text{WMMD}}^2(w) \right| + \left| \widehat{\text{WMMD}}^2(w) - \text{MMD}^2 \right| \end{aligned} \quad (40)$$

154 Then using Lemma 3 to bound the first term and Lemma 4 to bound the second term, we get that with  
155 probability at least  $1 - \delta$ ,

$$\begin{aligned} & \left| \widehat{\text{WMMD}}^2(\hat{w}) - \text{MMD}^2 \right| \leq \\ & R(B+1) \left( 2\epsilon + 2 \left( 1 + \sqrt{2 \log \frac{2}{\delta}} \right) R \sqrt{\frac{B^2}{n} + \frac{1}{m}} + (B+1) \sqrt{\frac{1}{2m_2} \log \frac{1}{\delta}} \right) \end{aligned} \quad (41)$$

156 where  $m_2 := \lfloor m/2 \rfloor$ .  $\square$

## 157 B. Details on the introductory example

158 The example is used to illustrate the need for adjusting for confounding variables. For a total of 500  
159 individuals we generated random education data  $X$  by sampling from a uniform distribution between  
160 0 and 10, from which we derived the post-intervention income  $Y^0$  and  $Y^1$  by simply adding a standard

161 random Gaussian noise variable to these values (in this case  $\mathcal{H}_0$  holds: the distributions are equal). We  
 162 generated male  $T = 1$  and female  $T = 0$  data, our two populations ( $S = 1$ ), by selectively removing  
 163 with probability 0.5 females with education level higher than 5 ( $Pr(T = 0|X > 5) \approx 0.33$ ), and  
 164 removing with probability 0.5 males with education level lower than 5 ( $Pr(T = 0|X < 5) \approx 0.66$ ).  
 165 We end up with approximately 150 individuals in each group, males with higher education levels than  
 166 females on average. Observe that the underlying generating process is the same in both populations,  
 167 only the marginal distribution of the education level changes. As is natural, a two-sample test that  
 168 overlooks the differences in education will reject the hypothesis of equal data generating process for  
 169 the income.

## 170 C. Description and implementation of tests

### 171 C.1. Hyperparameter selection for high power

172 The population quantity  $W\overline{MMD} = 0$  if and only if the distributions under consideration are equal,  
 173 for any choice of kernel hyperparameters. With finite sample size  $n$ , decisions must rely on inference  
 174 based on the empirical  $W\overline{MMD}$ , and some hyperparameters will give higher power than others. A  
 175 popular strategy is to set the bandwidth  $\sigma$  of the Gaussian kernel to the median squared pairwise  
 176 distance between input data, but can be sub-optimal when the scale of the difference between  
 177 populations differs from the scale of the difference within populations themselves. Instead, we follow  
 178 the approaches of [12, 7] and choose  $\sigma$  so as to maximize the test power, i.e. the probability of  
 179 rejecting  $\mathcal{H}_1$  when it is false.

180 **Proposition** (Approximate power of test statistic). *Under  $\mathcal{H}_1$ , for large  $n$  and fixed  $r$ , the test power*  
 181  $Pr(n\widehat{W\overline{MMD}}^2 > r) \approx 1 - \Phi\left(\frac{r}{\sqrt{n}\sigma_{\mathcal{H}_1}} - \sqrt{n}\frac{W\overline{MMD}}{\sigma_{\mathcal{H}_1}}\right)$ , *where  $\Phi$  denotes the cumulative distribution*  
 182 *function of the standard normal distribution, and  $\sigma_{\mathcal{H}_1}$  is defined as in Theorem 1.*

183 Assume that  $n$  is sufficiently large. Following the same argument as in [7], in  $\frac{r}{\sqrt{n}\sigma_{\mathcal{H}_1}} - \frac{W\overline{MMD}}{\sigma_{\mathcal{H}_1}}$ , we  
 184 observe that the first term  $\frac{r}{\sqrt{n}\sigma_{\mathcal{H}_1}} = \mathcal{O}(n^{-1/2})$  goes to 0 as  $n \rightarrow \infty$  because  $\sigma_{\mathcal{H}_1}^2 = \mathcal{O}(n^{-1})$ , while  
 185 the second term,  $\sqrt{n}\frac{W\overline{MMD}}{\sigma_{\mathcal{H}_1}} = \mathcal{O}(n^{1/2})$ , dominates the first one for large  $n$ . Thus, the parameters  
 186 that maximize the test power are given by  $\theta^* = \operatorname{argmax}_{\theta} p(n\widehat{W\overline{MMD}}^2 > r) \approx \frac{W\overline{MMD}}{\sigma_{\mathcal{H}_1}}$ . Since  
 187  $W\overline{MMD}$  and  $\sigma_{\mathcal{H}_1}$  are unknown, to maintain the validity of the hypothesis test we divide the sample  
 188 into a training set, used to compute  $\frac{\widehat{W\overline{MMD}}}{\hat{\sigma}_{\mathcal{H}_1}}$  and choose the kernel, and a testing set used to perform  
 189 the final hypothesis test with the learned kernel. The empirical estimate of the variance  $\hat{\sigma}_{\mathcal{H}_1}$  that  
 190 appears in our objective is approximated up to second order terms, similarly to [12].

### 191 C.2. B-Test: A modification that uses propensity scores

192 An alternative to the weighted MMD test is a B-test (block-based test): the idea is to break the  
 193 data into homogeneous blocks by stratifying subjects into mutually exclusive subsets based on their  
 194 estimated propensity score. Recall that the propensity score is defined as  $e(x) := Pr(T = 1|X)$ , the  
 195 probability of group assignment given confounding variables. After this stage, we compute a two  
 196 sample test statistic on each block, and average these quantities to obtain the test statistic.

197 More specifically, subjects are ranked according to their estimated propensity score and then stratified  
 198 into subsets based on previously defined thresholds of the estimated propensity score. Because  
 199 population assignment is essentially at random for individuals with the same propensity value, we  
 200 expect mean comparisons within this group to be unbiased. [9] showed that stratification based on the  
 201 propensity score will balance  $x$ , in the sense that within strata homogeneous in  $e(x) = Pr(T = 1|x)$ ,  
 202 the distribution of  $x$  will be equal in the two populations.

203 For an individual block, laying on the main diagonal and starting at position  $(i-1)B + 1$ , the statistic  
 204  $\eta(i)$  is calculated as,

$$\eta(i) := \frac{1}{\binom{B}{2}} \sum_{a=(i-1)B+1}^{iB} \sum_{b=(i-1)B+1 \neq a}^{iB} h(y_{a,0}, y_{b,0}^*, y_{a,1}, y_{b,1}^*) \quad (42)$$



205 where  $h(y_0, y_0^*, y_1, y_1^*) = k(y_0, y_0^*) + k(y_1, y_1^*) - k(y_0, y_1^*) - k(y_0^*, y_1)$ ,  $y_0$  is a sample from  
 206  $Y|T = 0$ ,  $y_1$  a sample from  $Y|T = 1$  and superscript  $\star$  denotes an independent copy. The overall  
 207 test statistic is then,

$$\eta = \frac{B}{n} \sum_{i=1}^{\frac{n}{B}} \eta(i) \quad (43)$$

208 The choice of  $B$  determines the accuracy of the balancing procedure and computation time - at one  
 209 extreme is exact matching based on the propensity score and the linear-time MMD suggested by [3]  
 210 where we have  $n/2$  blocks of size  $B = 2$ , and at the other extreme is the unbalanced and usual full  
 211 MMD with 1 block of size  $n$ . We chose as a default to divide both populations into  $\sqrt{n}$  blocks as  
 212 proposed in [15].

213 B-test of [15] assumes that  $B \rightarrow \infty$  together with  $n$ , which implies that the statistic  $\hat{\eta}$  defined in (43)  
 214 under the null distribution satisfies,

$$\sqrt{nB}\hat{\eta} \rightarrow_d \mathcal{N}(0, 4\sigma^2) \quad (44)$$

215 where  $\sigma^2 = E_{X, X'}(k(X, X')^2) + (E_{X, X'}k(X, X'))^2 - 2E_X[(E_{X'}k(X, X'))^2]$  that can be estimated  
 216 directly or by considering the empirical variance of the statistics computed within each of the blocks.

### 217 C.3. ANCOVA

218 Analysis of covariance (ANCOVA) are a general statistical procedure derived from a general linear  
 219 model which blend ANOVA and regression. Conventionally, ANCOVA evaluates whether the means  
 220 of a dependent variable are equal across levels of a categorical independent variable often called a  
 221 treatment, while statistically controlling for the effects of other continuous variables that are not of  
 222 primary interest, that is confounders. In existing implementations [13] these suffer from a number of  
 223 limitations such as the assumption of an underlying linear feature/outcome mapping and normality of  
 224 residuals.

225 In our implementation we proceed as follows. We fit a Random Forest regression model on the  
 226 confounding variables to approximate the outcome variable  $Y$ . Since in our experiments we consider  
 227  $Y$  to be multivariate, we fit a different regression model for each dimension of  $Y$ . We interpret the  
 228 resulting residuals as being independent of confounders given group assignments and use those to  
 229 proceed with testing. Because of the computational burden of this procedure, we fit the well-known  
 230 Hotelling  $T^2$  test [6] on the residuals to decide whether  $Y^0$  and  $Y^1$  share the same generating process  
 231 up to confounding variables.

## 232 D. Computational complexity

233 The computational complexity of the WMMD<sup>2</sup> is quadratic in the number of samples due to the  
 234 need to compute the Kernel matrix, similarly to the plain implementation of the MMD<sup>2</sup>. When  
 235 permutations are chosen to approximate the null distribution, this procedure can be overly time  
 236 consuming for large data sets. Below we briefly describe existing approximations that can be used  
 237 with the WMMD<sup>2</sup> to speed up computations.

- 238 • Gamma approximation to the null [4]. This procedure consist of using a two-parameter  
 239 Gamma distribution that we fit by matching the first and second moments of the empirical  
 240 MMD<sup>2</sup>. Such approximations can be accurate in practice and much faster, although they  
 241 remain heuristics with no consistency guarantees.
- 242 • Linear time test [3]. Another alternative would be to randomly subsample the data such as to  
 243 make the computational complexity linear in the original number of samples. The drawback  
 244 is that power is often overly reduced as a result.
- 245 • Kernel matrix approximation with low-dimensional random features [8]. To accelerate  
 246 the computation of the kernel matrix, one may map the input data to a randomized low-  
 247 dimensional feature space and compute inner products based on these representations. [8]  
 248 showed that by projecting unto a suitable basis the inner products of the transformed data are  
 249 approximately equal to those in the feature space of a user specified shift-invariant kernel.

## 250 References

- 251 [1] T De Wet, JH Venter, et al. Asymptotic distributions for quadratic forms with applications to  
252 tests of fit. *The Annals of Statistics*, 1(2):380–387, 1973.
- 253 [2] Arthur Gretton, Karsten Borgwardt, Malte Rasch, Bernhard Schölkopf, and Alex J Smola. A  
254 kernel method for the two-sample-problem. In *Advances in neural information processing  
255 systems*, pages 513–520, 2007.
- 256 [3] Arthur Gretton, Karsten M Borgwardt, Malte J Rasch, Bernhard Schölkopf, and Alexander  
257 Smola. A kernel two-sample test. *Journal of Machine Learning Research*, 13(Mar):723–773,  
258 2012.
- 259 [4] Arthur Gretton, Kenji Fukumizu, Zaid Harchaoui, and Bharath K Sriperumbudur. A fast,  
260 consistent kernel two-sample test. In *Advances in neural information processing systems*, pages  
261 673–681, 2009.
- 262 [5] Arthur Gretton, Alex Smola, Jiayuan Huang, Marcel Schmittfull, Karsten Borgwardt, and  
263 Bernhard Schölkopf. Covariate shift by kernel mean matching. *Dataset shift in machine  
264 learning*, 2009.
- 265 [6] Harold Hotelling. The generalization of student’s ratio. In *Breakthroughs in statistics*, pages  
266 54–65. Springer, 1992.
- 267 [7] Wittawat Jitkrittum, Wenkai Xu, Zoltán Szabó, Kenji Fukumizu, and Arthur Gretton. A linear-  
268 time kernel goodness-of-fit test. In *Advances in Neural Information Processing Systems*, pages  
269 262–271, 2017.
- 270 [8] Ali Rahimi and Benjamin Recht. Random features for large-scale kernel machines. In *Advances  
271 in neural information processing systems*, pages 1177–1184, 2008.
- 272 [9] Paul R Rosenbaum and Donald B Rubin. The central role of the propensity score in observational  
273 studies for causal effects. *Biometrika*, 70(1):41–55, 1983.
- 274 [10] Robert J Serfling. *Approximation theorems of mathematical statistics*, volume 162. John Wiley  
275 & Sons, 2009.
- 276 [11] Connie P Shapiro, Lawrence Hubert, et al. Asymptotic normality of permutation statistics  
277 derived from weighted sums of bivariate functions. *The Annals of Statistics*, 7(4):788–794,  
278 1979.
- 279 [12] Dougal J Sutherland, Hsiao-Yu Tung, Heiko Strathmann, Soumyajit De, Aaditya Ramdas, Alex  
280 Smola, and Arthur Gretton. Generative models and model criticism via optimized maximum  
281 mean discrepancy. *arXiv preprint arXiv:1611.04488*, 2016.
- 282 [13] Barbara G Tabachnick, Linda S Fidell, and Jodie B Ullman. *Using multivariate statistics*,  
283 volume 7. Pearson Boston, MA, 2019.
- 284 [14] Steve Verrill and Richard A Johnson. Asymptotic distributions for quadratic forms with  
285 applications to censored data tests of fit. *Communications in Statistics-Theory and Methods*,  
286 17(12):4011–4024, 1988.
- 287 [15] Wojciech Zaremba, Arthur Gretton, and Matthew Blaschko. B-test: A non-parametric, low  
288 variance kernel two-sample test. In *Advances in neural information processing systems*, pages  
289 755–763, 2013.