

---

# AutoML Two-Sample Test

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 Two-sample tests are important in statistics and machine learning, both as tools  
2 for scientific discovery as well as to detect distribution shifts. This led to the  
3 development of many sophisticated test procedures going beyond the standard  
4 supervised learning frameworks, whose usage can require specialized knowledge  
5 about two-sample testing. We use a simple test that takes the mean discrepancy  
6 of a witness function as the test statistic and prove that minimizing a squared loss  
7 leads to a witness with optimal testing power. This allows us to leverage recent  
8 advancements in AutoML. Without any user input about the problems at hand,  
9 and using the same method for all our experiments, our AutoML two-sample test  
10 achieves competitive performance on a diverse distribution shift benchmark as well  
11 as on challenging two-sample testing problems.

## 12 1 Introduction

13 Testing whether two distributions are the same based on data is a fundamental problem in data science.  
14 A classical application is to test whether two differently treated groups have the same characteristics  
15 or not [Student, 1908, Welch, 1947, Golland and Fischl, 2003]. Testing independence of two random  
16 variables can also be phrased as a two-sample problem by testing whether the joint distribution equals  
17 the product of the marginal distributions [Gretton et al., 2005]. A more recent application in machine  
18 learning is to detect distribution shifts, i.e., whether the distribution a model was trained on equals the  
19 distribution the model is deployed on [Lipton et al., 2018, Rabanser et al., 2019, Koch et al., 2022].

20 Classical methods have a fixed test statistic that makes strong parametric assumptions. For example,  
21 Student’s two-sample  $t$ -test only tests whether the distributions have equal mean, assuming both  
22 distributions follow a normal distribution with the same (but unknown) variance. With modern  
23 datasets, which are often high-dimensional, such test cannot be applied because the strong assumptions  
24 are often not justified. Nonparametric kernel-based test such as the Maximum Mean Discrepancy  
25 (MMD) [Gretton et al., 2012a] are very flexible and, theoretically, can detect differences of any kind  
26 given enough data. However, this generality often harms test power at finite data size. This can  
27 simply be understood in terms of a classical bias-variance tradeoff. To mitigate this, it has become  
28 common to optimize the test statistic first on a held-out dataset and then apply the test only on the  
29 other part of the data [Sutherland et al., 2017, Liu et al., 2020]. However, the derived objective as  
30 well as optimizing a kernel function are no standard tasks in machine learning and no automated  
31 packages exist, making it hard for practitioners to apply them.

32 Tests that fit well into the standard machine learning pipeline are based on the classification accuracy.  
33 First, a classifier is trained to detect the difference between the two samples, and then its accuracy on  
34 a held-out set is taken as a test statistic [Golland and Fischl, 2003, Lopez-Paz and Oquab, 2017, Kim  
35 et al., 2021, Cai et al., 2020, Hediger et al., 2022]. Liu et al. [2020] argued, however, that optimizing  
36 classification accuracy does not directly optimize test power and considered this one reason why  
37 kernel-based test outperform classifier tests. Kübler et al. [2022] challenged this and considered the  
38 mean of an optimized witness function as test statistic finding that kernels are not necessary for good

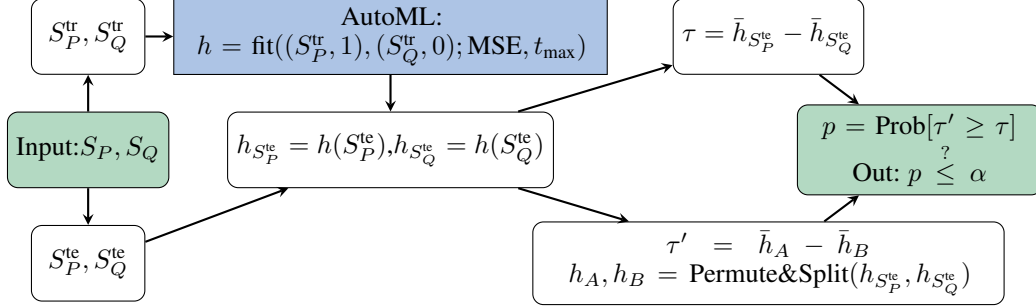


Figure 1: AutoML two-sample test:  $S_P, S_Q$  denotes the available data from  $P$  and  $Q$ , which is first split into two parts of equal size. A witness  $h : \mathcal{X} \rightarrow \mathbb{R}$  is trained using a (weighted) squared loss Eq. (6), denoted by MSE, and using AutoML to maximize predictive performance. Users can easily control important properties, for example the maximal runtime  $t_{\max}$ . The test statistic  $\tau$  is the difference in means on the test sets. Permuting the data and recomputing  $\tau$  allows the estimation of the  $p$ -values. The null hypothesis  $P = Q$  is rejected if  $p \leq \alpha$ .

39 performance. Generally, such two-stage procedures are very intuitive and arguably also how a human  
 40 would approach the two-sample problem on complicated data. One could look at some part of the  
 41 data, try to come up with a simple hypothesis, and then try to test its significance on held-out data.

42 Despite the recent progress in the theoretical understanding of machine learning-based two sample  
 43 tests [Kim et al., 2021, Liu et al., 2020], there is still little guidance on how to apply these tests in  
 44 practice and a substantial amount of engineering and expertise is required to implement them. On  
 45 the contrary, in supervised learning, namely regression and classification, the past years have shown  
 46 tremendous advancements in making machine learning models applicable essentially without any  
 47 expert knowledge leading to the field of Automated Machine Learning (AutoML) [Feurer et al., 2015,  
 48 Hutter et al., 2019, He et al., 2021]. The goal of AutoML is to automate the full machine learning  
 49 pipeline: Data cleaning, feature engineering and augmentation, model search, hyperparameter  
 50 optimization, and model ensembling [Dietterich, 2000]. All of it with the goal of achieving the best  
 51 possible predictive performance on unseen data.

52 The goal of our work is to bring the advancements of AutoML research to the field of two-sample  
 53 testing. Our main contributions are:

- 54 1. We prove that minimizing a squared loss is equivalent to maximizing the unwieldy signal-to-noise  
 55 ratio, which determines the asymptotic test power of a witness two-sample test (Section 3.1).
- 56 2. Thanks to the former result we can use AutoML to learn the test statistic, thereby harnessing the  
 57 power of many advancements in machine learning such as hyperparameter optimization, bagging,  
 58 and ensemble learning in a user-friendly manner (Section 3.2).
- 59 3. Our test is usable without any specific knowledge and skills in two-sample testing. Users can easily  
 60 specify how many resources they want to use when learning the test, for example the maximal training  
 61 time (Section 3.2 and Section 5). Furthermore, one can easily interpret the results (Section 3.3).
- 62 4. We extensively study the empirical performance of our approach first by considering the two  
 63 low-dimensional datasets Blob and Higgs followed by running a large benchmark on a variety of  
 64 distribution shifts on MNIST and CIFAR10 data. We observe very competitive performance without  
 65 any manual adjustment of hyperparameters. Our experiments also show that a continuous witness  
 66 outperforms commonly used binary classifiers (Section 5).

67 The proposed testing pipeline is described in Fig. 1: First, the two samples are split into training and  
 68 test sets. Then a *witness* function  $h$  is trained by first labeling samples in  $S_P^{\text{tr}}$  with 1 and samples from  
 69  $S_Q^{\text{tr}}$  with 0 and then minimizing a (weighted) Mean Squared Error (MSE) to maximize test power,  
 70 see Section 3 for further details. To maximize the predictive performance and to require as little  
 71 user input as possible, we use AutoGluon [Erickson et al., 2020], an existing AutoML framework,  
 72 when optimizing the witness. Our test statistic is then simply the difference in means of the test sets  
 73  $S_P^{\text{te}}, S_Q^{\text{te}}$ , see Section 2.  $p$ -values are computed via permutation of the samples [Golland and Fischl,  
 74 2003], which is a standard technique in two-sample testing.

75 **2 Preliminaries**

76 **Notation.** We consider the non-empty set  $\mathcal{X} \subseteq \mathbb{R}^d$  as the domain of our data. We are given samples  
 77  $S_P = \{x_1, \dots, x_n\}$  and  $S_Q = \{y_1, \dots, y_m\}$ , which are i.i.d. realizations of the random variables  
 78  $X$  and  $Y$  distributed according to  $P$  and  $Q$ . Let us define  $c = \frac{n}{n+m}$ . The proposed approach splits  
 79 the data into disjoint training and test sets of size  $n_{tr}, n_{te}, m_{tr}, m_{te}$  that we denote  $S_P^{tr}, S_P^{te}$  and  $S_Q^{tr}, S_Q^{te}$ .  
 80 Unless otherwise stated, we assume that the data is split in equal halves, which is the default approach  
 81 [Lopez-Paz and Oquab, 2017, Liu et al., 2020].

82 Our goal is to test the null hypothesis  $H_0 : P = Q$  against the alternative hypothesis  $P \neq Q$ . Our  
 83 hypothesis test rejects when the observed value of the test statistic is significantly larger than what  
 84 we would expect if the null hypothesis were true. Naturally, such tests can make two types of errors:

85 Type-I: The test rejects the null hypothesis, although it is true.

86 Type-II: The test fails to reject, although the null hypothesis is false.

87 Our goal is to design a test that controls the rate of Type-I errors at a given significance level  
 88  $\alpha \in (0, 1)$ , and maximizes the test *power* defined as 1 minus the probability of a Type-II error.

89 **Witness two-sample test.** We consider a witness-based hypothesis test [Cheng and Cloninger,  
 90 2019, Kübler et al., 2022]. Given a function  $h : \mathcal{X} \rightarrow \mathbb{R}$ , called *witness*, the *mean discrepancy* is

$$\tau(P, Q | h) = \mathbb{E}_{X \sim P} [h(X)] - \mathbb{E}_{Y \sim Q} [h(Y)], \quad (1)$$

91 and we use its empirical estimate on the test set as test statistic

$$\tau(S_P^{te}, S_Q^{te} | h) = \frac{1}{n_{te}} \sum_{x \in S_P^{te}} h(x) - \frac{1}{m_{te}} \sum_{y \in S_Q^{te}} h(y). \quad (2)$$

92 As we show in Section 4, this test statistic can be seen as a continuous extension of classifier two-  
 93 sample tests [Lopez-Paz and Oquab, 2017]. We assume that  $c = \frac{n_{te}}{n_{te} + m_{te}}$  converges to a constant. With  
 94  $\sigma_c^2(h) = \frac{(1-c)\text{Var}_{X \sim P}[h(X)] + c\text{Var}_{Y \sim Q}[h(Y)]}{c(1-c)}$  we have that the test statistic is asymptotically normally  
 95 distributed [Kübler et al., 2022, Theorem 1]

$$\sqrt{n_{te} + m_{te}} [\tau(S_P^{te}, S_Q^{te} | h) - \tau(P, Q | h)] \xrightarrow{d} \mathcal{N}(0, \sigma_c^2(h)). \quad (3)$$

96 Let us for now assume that we know  $\sigma_c(h)$ . For any level  $\alpha \in (0, 1)$  we can set the *analytic* test  
 97 threshold to  $t_\alpha = \frac{\sigma_c(h)}{\sqrt{n_{te} + m_{te}}} \Phi^{-1}(1 - \alpha)$ , where  $\Phi$  denotes the CDF of a standard normal and  $\Phi^{-1}$  its  
 98 inverse. We can then compute the asymptotic probability of rejecting as:

$$\Pr[\text{reject}] = \Pr[\tau(S_P^{te}, S_Q^{te} | h) > t_\alpha] \rightarrow \Phi\left(\frac{\sqrt{n_{te} + m_{te}} \tau(P, Q | h)}{\sigma_c(h)} - \Phi^{-1}(1 - \alpha)\right). \quad (4)$$

99 Under the null hypothesis  $P = Q$  we have  $\tau(P, Q | h) = 0$ . Therefore, Eq. (4) reduces to  
 100  $\Phi(-\Phi^{-1}(1 - \alpha)) = 1 - \Phi(\Phi^{-1}(1 - \alpha)) = \alpha$ . Hence, the asymptotic test correctly controls type-I  
 101 error. On the other hand, for given  $P \neq Q$ , Eq. (4) corresponds to the test power. Since  $\Phi$  is a  
 102 monotonically increasing function, the test power is maximized by the witness  $h$  that maximizes

$$\text{SNR}(h) = \frac{\tau(P, Q | h)}{\sigma_c(h)}, \quad (5)$$

103 where SNR is the Signal-to-Noise Ratio. Kübler et al. [2022] showed the optimal witness can be  
 104 learned when using kernel methods (using kernel Fisher Discriminant Analysis), but it was left  
 105 open how this can be done efficiently with other machine learning frameworks. Such an SNR is not  
 106 commonly implemented and also common approaches like mini-batching are not easily adapted, as  
 107 an estimate of the SNR based on a mini batch would be a biased estimate. In the next section, we  
 108 show how to circumvent this and optimize a squared loss instead.

### 109 3 The AutoML two-sample test

#### 110 3.1 Equivalence of squared loss and signal-noise ratio

111 Since it is known for linear models that minimizing a squared loss over two labelled samples is  
 112 equivalent to Fisher Discriminant Analysis [Duda et al., 2001, Mika, 2003], we attempt to find a  
 113 more general relation between the squared loss and the SNR. Our goal is to use the squared loss as  
 114 the optimization objective when learning the witness. Let  $c = \frac{n_u}{n_u + m_u}$  analogously to the above. Let  
 115 us mark all data from  $P$  with a label '1' and all data from  $Q$  with a label '0'. We define the following  
 116 (weighted) squared loss

$$L_{P,Q,c}(h) = (1 - c) \mathbb{E}_{X \sim P} [(1 - h(X))^2] + c \mathbb{E}_{Y \sim Q} [(0 - h(Y))^2], \quad (6)$$

117 Note that the weights  $(1 - c)$  and  $c$  are swapped as it will be more important to fit the set with fewer  
 118 samples. Given a function  $h$ , notice that shifting and scaling it leaves the SNR (5) invariant. We can  
 119 then show the following relationship of its squared loss and its SNR.

120 **Lemma 1.** *Let the function  $h$  be fixed. We apply the linear transformation  $h \rightarrow \gamma h + \nu$  with  $\gamma \in \mathbb{R}$   
 121 and  $\nu \in \mathbb{R}$ . Let  $(\gamma^*, \nu^*)$  be the minimum of the quadratic function  $(\gamma, \nu) \mapsto L(\gamma h + \nu)$ . Then, the  
 122 following holds true:*

$$L(\gamma^* h + \nu^*) = \frac{c(1 - c)}{1 + \text{SNR}(h)^2}.$$

123 *Proof.* We refer the reader to Appendix A. □

124 Let us assume that the supports of the two distributions  $P, Q$  overlap. Hence, for any function the loss  
 125  $L_{P,Q,c}$  is strictly positive. Assume that  $h^*$  is the function that minimizes the loss over all possible  
 126 functions. This implies that  $\gamma^* = 1$  and  $\nu^* = 0$ , as otherwise one could still improve the loss by  
 127 scaling or shifting. Thus, by Lemma 1 we have:

128 **Proposition 1.** *Assume that  $h^*$  minimizes the squared loss (6). Then  $h^*$  maximizes the signal-to-noise  
 129 ratio, i.e.,*

$$L(h^*) = \min_h L(h) \Rightarrow \text{SNR}(h^*) = \max_h \text{SNR}(h).$$

130 *Proof.* A solution that minimizes the loss has  $\bar{h}_P^* \geq \bar{h}_Q^*$  and hence a non-negative SNR. Assume  
 131 there exists  $\tilde{h}$  such that  $\text{SNR}(\tilde{h}) > \text{SNR}(h^*)$ . Then Lemma 1 implies the existence of  $\tilde{\gamma}, \tilde{\nu}$  such that  
 132  $L(\tilde{\gamma}\tilde{h} + \tilde{\nu}) < L(h^*)$ , which is a contradiction. □

133 We can further derive a closed-form solution for the population optimal witness:

134 **Proposition 2 (Optimal Witness).** *Assume  $P$  and  $Q$  have densities  $p(x)$  and  $q(x)$ . The function  
 135 minimizing Eq. (6) is*

$$h^*(x) = \frac{(1 - c)p(x)}{(1 - c)p(x) + cq(x)}. \quad (7)$$

136 *Proof.* We rewrite Eq. (6) as  $L(h) = \int_{\mathcal{X}} (1 - c)p(x)(1 - h(x))^2 + cq(x)h^2(x) dx$ . Minimizing the  
 137 integrand for each  $x$  yields the claimed result. A similar result was obtained by Mao et al. [2019]. □

138 **Remark 1.** *Consider the balanced case  $c = 1/2$ , i.e., equal prior probabilities of labels '1' and '0'.  
 139 Then  $h^*(x)$  is the posterior probability that the example  $x$  came from  $P$ , or, using our defined labels,  
 140  $h^*(x) = \text{Pr}[1|x]$ . Thus, minimizing a log loss, i.e. the binary cross-entropy, and using its output  
 141 probability for class 1 as witness function also maximizes test power.*

142 Notice that for  $c \neq 1/2$ , we need to weight our samples with the inverse weights, i.e., it is more  
 143 important to get the less frequent samples right.

144 Proposition 1 and Remark 1 lead to our main conclusion: *To find an optimal witness, we can simply  
 145 optimize the (weighted) squared error or a cross-entropy loss.* This allows us to seamlessly integrate  
 146 existing AutoML frameworks, which are designed to solve this task in an automated fashion, to learn  
 147 powerful witnesses. In the following we mainly focus on the squared error.

148 **3.2 Practical implementation**

149 **Stage I - optimization.** In the first stage, we optimize the witness function to minimize the MSE  
 150 via the training data  $S_P^{\text{tr}}$  and  $S_Q^{\text{tr}}$ , as motivated in the previous section. We simply label the data with  
 151 1 or 0 depending on whether they come from  $P$  or  $Q$ . We can then use any library that implements  
 152 an optimization of a squared loss. If  $c \neq 1/2$  we additionally need to specify weights according  
 153 to Eq. (6). Note that, unsurprisingly, the relevant quantity for the test power is the loss on the test  
 154 data and not on the training data. Thus, it is of crucial importance to find a witness with good  
 155 generalization performance. To make this as simple as possible for practitioners, we propose to use  
 156 an AutoML framework. This also has the advantage that users can specify runtime and memory  
 157 limits, and can explicitly trade computational resources for better statistical significance.

158 Although we strongly argue towards using AutoML for the test, this can of course not circumvent  
 159 the no-free-lunch theorem. Thus, whenever users have good intuition about how their two samples  
 160 might differ, we strongly encourage taking this into account when designing the test. To put it to  
 161 the extreme: If one knows that their (one-dimensional) data follows a normal distribution and only  
 162 differs in mean (if at all), one should use a classic  $t$ -test rather than our approach.

163 **Stage II - testing.** Given a witness function  $h$  learned as detailed in the previous section, we compute  
 164 the test statistic as in Eq. (2). To compute a p-value or decide whether to reject the null hypothesis  $P =$   
 165  $Q$ , we can either approximate the asymptotic distribution or use permutations [Kübler et al., 2022].  
 166 To estimate an asymptotically valid p-value<sup>1</sup> we first estimate  $\sigma_c^2(h)$  (see Eq. (3)) based on  $S_P^{\text{te}}, S_Q^{\text{te}}$ ,  
 167 which we denote as  $\hat{\sigma}_c^2(h)$ . The p value is then given as  $1 - \Phi\left(\sqrt{n_{\text{te}} + m_{\text{te}}}\tau(S_P^{\text{te}}, S_Q^{\text{te}}|h)/\hat{\sigma}_c(h)\right)$ .

168 For two-sample tests, a cheap alternative that guarantees correct type-I error control even at finite  
 169 sample size is based on permutations [Golland and Fischl, 2003]. In case of witness functions, one  
 170 can simply permute the values  $h(x_1), \dots, h(x_{n_{\text{te}}}), h(y_1), \dots, h(y_{m_{\text{te}}})$  and split them in two sets of  
 171 size  $n_{\text{te}}$  and  $m_{\text{te}}$ , respectively. One then recomputes the test statistic and estimates the  $p$ -value based  
 172 on the empirical quantile over  $B \in \mathbb{N}$  iterations. We reject whenever  $p \leq \alpha$ . We emphasize that we  
 173 do not need to retrain the model, and it even suffices to evaluate the witness once on all elements of  
 174 the test sets. We can then directly permute the witness' values.

175 **Runtime.** The overall runtime of the AutoML based witness test is the sum of the runtimes of  
 176 the training phase, the evaluation of the witness, and the evaluation of the test statistic. We denote  
 177 the scaling of the former by  $s_{\text{train}}[n_{\text{tr}} + m_{\text{tr}}]$ , where the square brackets emphasize that this is not a  
 178 product. It will depend on the AutoML framework but can usually be controlled by setting a time  
 179 limit. Even with a limit of one minute or less AutoGluon can already train powerful models on large  
 180 datasets and even performs model-selection, hyperparameter optimization, and so on. In contrast,  
 181 deep kernel-based methods typically train a neural network with a fixed architecture, which can  
 182 be expensive. Although neural networks belong to the suite of models AutoGluon trains, they are  
 183 optimized for speed and if the runtime limit does not permit training them another faster model will  
 184 be selected.

185 The scaling of evaluating  $h$ , denoted by  $s_{\text{eval}}[n_{\text{tr}} + m_{\text{tr}}]$ , is usually linear in the dataset size, but  
 186 it can be sublinear if the evaluation is parallelized. It can also be controlled with AutoGluon  
 187 by using different hyperparameter presets which might optimize the model selection towards fast  
 188 inference times. Compared to that, kernel-based tests have a quadratic runtime. Furthermore, the test  
 189 statistic has to be evaluated on the original partition of the data as well as  $B$  permutations requiring  
 190  $(n_{\text{tr}} + m_{\text{tr}})(B + 1)$  steps. In practice, this is usually the cheapest step, but it could also be further  
 191 reduced by parallelization. The overall runtime is given by

$$O(s_{\text{train}}[n_{\text{tr}} + m_{\text{tr}}] + s_{\text{eval}}[n_{\text{te}} + m_{\text{te}}] + (n_{\text{te}} + m_{\text{te}})(B + 1)). \quad (8)$$

192 Generally, training the witness will be the most expensive step of our test. A main advantage of  
 193 our test over others is that practitioners can easily trade-off spending more time and resources on  
 194 the training phase to potentially get a better witness and thus to more significant results. Thanks to  
 195 AutoML, specifying the time and resources does not require any detailed knowledge of the underlying  
 196 algorithm and is hence easily done.

<sup>1</sup>Asymptotic p-values are strictly speaking only valid for fixed  $h$  as the size of  $S_P^{\text{te}}, S_Q^{\text{te}}$  goes to infinity.



Figure 2: Testing MNIST against shifted MNIST with '0's knocked out. The optimized witness assigns the highest values to the images on the left, and lowest values to the images on the right, allowing us to interpret the difference.

197 **3.3 Interpretability**

198 Suppose our test finds a significant difference between  $S_P$  and  $S_Q$ . An additional task would be  
 199 to *interpret* how the distributions differ. This is particularly simple in our framework and shown in  
 200 Fig. 2: We can check which examples attained the highest value of the witness to find which inputs  
 201 are much more likely under  $P$  than under  $Q$ . On the other hand, inputs with small witness values are  
 202 more likely under  $Q$ . Similar procedures were used in Jitkrittum et al. [2016], Lopez-Paz and Oquab  
 203 [2017], Rabanser et al. [2019]. An additional advantage of using the AutoML framework AutoGluon  
 204 is that it allows to compute feature importance values easily. Therefore, for datasets which are hard  
 205 to visualize the important features of data points with high or low witness values can be identified.

206 **4 Relation to prior work**

207 Gretton et al. [2012a] introduced the maximum mean discrepancy as a test statistic for two-sample  
 208 testing. For a given reproducing kernel Hilbert space  $\mathcal{H}$ , the maximum mean discrepancy is

$$\text{MMD} = \max_{h \in \mathcal{H}, \|h\| \leq 1} \mathbb{E}_{X \sim P} [h(X)] - \mathbb{E}_{Y \sim Q} [h(Y)]. \quad (9)$$

209 Commonly an empirical estimate of the squared MMD is taken as test statistic and thresholds  
 210 estimated via permutations. The connection to our test is quite apparent, since both use the mean  
 211 discrepancy as test statistic - solely that the MMD optimizes the function over the RKHS unit ball,  
 212 and the witness test tries to maximize the test power when learning the function on a held-out data  
 213 set. The past years have shown that learning the kernel in a data-driven manner improves testing  
 214 performance [Gretton et al., 2012b, Sutherland et al., 2017, Liu et al., 2020, Kübler et al., 2020,  
 215 Liu et al., 2021]. Liu et al. [2020] proposed to learn a deep kernel, by deriving an asymptotic test  
 216 power criterion and considering a rather involved kernel function (MMD-D), see their Eq. (1). They  
 217 concluded that this is often better than learning a one-dimensional representation, like a classifier  
 218 does. Kübler et al. [2022] challenged this, by showing that learning a simple witness function via  
 219 a test power criterion often suffices. They showed how to use cross-validation and kernel Fisher  
 220 Discriminant Analysis (kfda) to find powerful witness functions (kfda-witness), which serves as the  
 221 blueprint for our more general approach. Chwialkowski et al. [2015] proposed another kernel-based  
 222 fast two-sample test with smooth characteristic functions (SCF) and projected mean embeddings  
 223 (ME), which was refined by Jitkrittum et al. [2016] who optimized this test statistic in the first stage.  
 224 Kirchler et al. [2020] trained a deep multidimensional representation and used its mean distance as  
 225 test statistic. Recently, Zhao et al. [2022] proposed a general framework that also includes MMD.

226 Classifier two-sample tests (C2ST) also rely on a data splitting approach and have extensively been  
 227 studied in the literature [Friedman, 2003, Golland and Fischl, 2003, Lopez-Paz and Oquab, 2017,  
 228 Kim et al., 2021, Cai et al., 2020, Hediger et al., 2022]. For simplicity, we focus on the balanced case.  
 229 A C2ST trains a classifier with  $S_P^{\text{tr}}$ , labelled with '1' and  $S_Q^{\text{tr}}$  labelled with '0' and then estimates its  
 230 classification accuracy on  $S_P^{\text{te}}, S_Q^{\text{te}}$ . If the estimated accuracy is significantly above chance (that's  
 231 what it would be under the null hypothesis), the test rejects. Let  $f : \mathcal{X} \rightarrow \{0, 1\}$  denote the binary  
 232 classifier, then we can write the accuracy as  $\frac{1}{2} + \varepsilon$  and estimate it as

$$\begin{aligned} \frac{1}{2} + \hat{\varepsilon} &= \frac{1}{2} \left( \frac{1}{n_{\text{te}}} \sum_{i=1}^{n_{\text{te}}} f(x_i) + \frac{1}{n_{\text{te}}} \sum_{i=1}^{n_{\text{te}}} (1 - f(y_i)) \right) = \frac{1}{2} + \frac{1}{2} \left( \frac{1}{n_{\text{te}}} \sum_{i=1}^{n_{\text{te}}} f(x_i) - \frac{1}{n_{\text{te}}} \sum_{i=1}^{n_{\text{te}}} f(y_i) \right) \\ &= \frac{1}{2} + \frac{1}{2} \tau(S_P^{\text{te}}, S_Q^{\text{te}} | f). \end{aligned}$$

233 Thus using the classification accuracy as test statistic is equivalent to using the mean discrepancy  
 234 as test statistic with the binary classifier  $f$  as witness function in Eq. (2). However, using binary

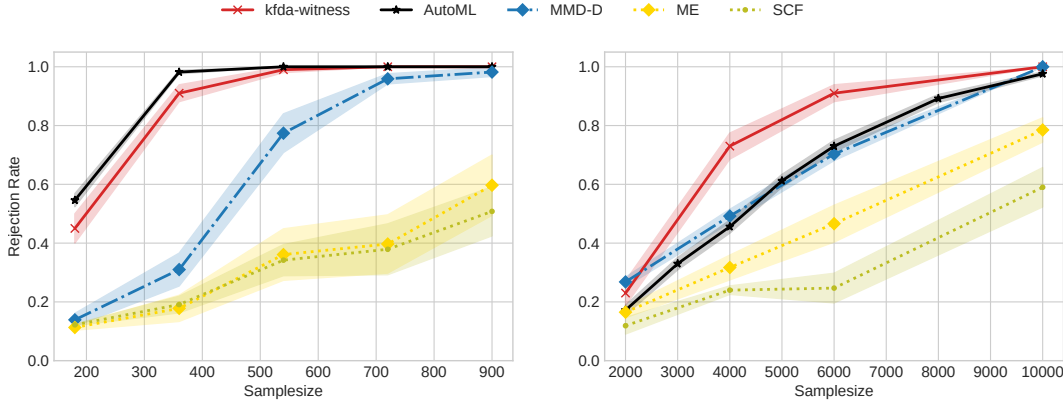


Figure 3: Experiments on low dimensional problems. The simple approach of learning a one-dimensional witness function with AutoML or optimizing a witness via kfda and a grid search can outperform more involved approaches. **Left: Blob, Right: Higgs.**

235 classifiers is quite limiting and results in quite high variance. Using continuous witness functions  
 236 allows for higher power.<sup>2</sup> Some might also speak of 'classifier' test when referring to a witness test,  
 237 but using the term 'witness' emphasizes that it is continuous.

238 Lipton et al. [2018] proposed to use a pretrained classifier to detect label shift. Rabanser et al. [2019]  
 239 extended this to detect covariate shift. They investigate different ways of reducing the dimensionality  
 240 and then applying different (classical) hypothesis test on them. While they also consider a basic C2ST,  
 241 their best performing method uses the softmax outputs of a pretrained image classifier. They then run  
 242 a *univariate* Kolmogorov-Smirnov test on each of the output 'probabilities' separately and correcting  
 243 via Bonferroni correction. We refer to this as (univariate) BBSDs (black box shift detection - soft).  
 244 For more details on their other methods, we refer the reader to their work directly.

## 245 5 Experiments

246 To show the power of utilizing AutoML we use the same setup for all datasets we consider. The  
 247 data is split into two equally sized parts since this is the standard approach [Lopez-Paz and Oquab,  
 248 2017, Liu et al., 2020, Rabanser et al., 2019]. We label data from  $P$  with '1', data from  $Q$  with '0'  
 249 and fit a least square regression with AutoGluon's TabularPredictor [Erickson et al., 2020]. We  
 250 use the configuration `presets='best_quality'` and by default optimize with a five-minute time  
 251 limit. For more details, we refer to the AutoGluon documentation. We run all experiments with  
 252 significance level  $\alpha = 5\%$ . Results of correct type-I error control are provided in Appendix B. The  
 253 sample size we report is always the size of the datasets before splitting, i.e.  $n = m$ , since we only  
 254 consider balanced problems.

255 All experiments were done on servers having only CPUs and we spend around 100k CPU hours on  
 256 doing all the experiments reported in the paper, which is mainly because we did various configurations  
 257 and many repetitions for all the test cases we consider. Further details are given in Appendix B.

258 **Blob & Higgs.** We first compare the performance on two low-dimensional datasets used to examine  
 259 the power of two-sample tests. Following Liu et al. [2020] we consider the Blob dataset, which is a  
 260 mixture of nine Gaussian modes, with different covariance structures between  $P$  and  $Q$  [Liu et al.,  
 261 2020, Fig. 1]. We also consider the Higgs dataset [Baldi et al., 2014], which was introduced by  
 262 Chwialkowski et al. [2015] as a two-sample problem. As baselines, we use MMD-D, ME, SCF, and  
 263 kfda-witness as reported by Kübler et al. [2022]. We report the results in Fig. 3, where  $\pm 1$  standard  
 264 error are shown as shaded regions. Since we estimated the performance over 500 runs, we obtain a  
 265 smaller error than the other methods. We observe that both approaches based on the mean difference  
 266 of a witness function (kfda-witness, AutoML) perform competitively. AutoML performs best on  
 267 Blob, and kfda-witness is best on Higgs.

<sup>2</sup>Lopez-Paz and Oquab [2017] also observe that using a binary classifier might be too restrictive (see their Remark 2), but they did not investigate this in detail.

Table 1: Shift detection on MNIST and CIFAR10 based on Rabanser et al. [2019].

(a) Test power across all simulated shifts on MNIST and CIFAR10. We propose the methods in the bottom section.

Test	DR	Number of samples from test							
		10	20	50	100	200	500	1,000	10,000
Univ. tests	NoRed	0.03	0.15	0.26	0.36	0.41	0.47	0.54	0.72
	PCA	0.11	0.15	0.30	0.36	0.41	0.46	0.54	0.63
	SRP	0.15	0.15	0.23	0.27	0.34	0.42	0.55	0.68
	UAE	0.12	0.16	0.27	0.33	0.41	0.49	0.56	0.77
	TAE	0.18	0.23	0.31	0.38	0.43	0.47	0.55	0.69
	BBSDs	0.19	0.28	<b>0.47</b>	<b>0.47</b>	0.51	<b>0.65</b>	<b>0.70</b>	0.79
$\chi^2$	BBSDh	0.03	0.07	0.12	0.22	0.22	0.40	0.46	0.57
	Classif	0.01	0.03	0.11	0.21	0.28	0.42	0.51	0.67
Multiv. tests	NoRed	0.14	0.15	0.22	0.28	0.32	0.44	0.55	-
	PCA	0.15	0.18	0.33	0.38	0.40	0.46	0.55	-
	SRP	0.12	0.18	0.23	0.31	0.31	0.44	0.54	-
	UAE	<b>0.20</b>	0.27	0.40	0.43	0.45	0.53	0.61	-
	TAE	0.18	0.26	0.37	0.38	0.45	0.52	0.59	-
	BBSDs	0.16	0.20	0.25	0.35	0.35	0.47	0.50	-
AutoGluon (raw)		0.17	0.24	0.37	0.46	0.50	0.62	0.67	<b>0.87</b>
AutoML (pre)		0.18	<b>0.29</b>	<b>0.42</b>	<b>0.47</b>	0.47	0.64	0.65	0.72
AutoML (class)		0.19	0.19	0.38	0.46	<b>0.52</b>	0.61	0.67	<b>0.87</b>
AutoML (bin)		0.03	0.14	0.31	0.43	0.49	0.51	0.59	0.86

Shift	Test	Number of samples from test							
		10	20	50	100	200	500	1,000	10,000
s_gn	raw	0.20	0.27	0.33	0.40	0.43	0.50	0.63	0.80
	pre	0.00	0.03	0.10	0.03	0.00	0.10	0.03	0.03
m_gn	raw	0.27	0.23	0.33	0.43	0.43	0.53	0.63	0.83
	pre	0.00	0.03	0.17	0.00	0.00	0.13	0.07	0.13
l_gn	raw	0.23	0.33	0.53	0.67	0.70	0.77	1.00	1.00
	pre	0.17	0.27	0.50	0.57	0.60	0.73	0.80	0.90
s_img	raw	0.13	0.27	0.30	0.33	0.40	0.50	0.53	0.83
	pre	0.20	0.30	0.60	0.57	0.67	0.83	0.83	1.00
m_img	raw	0.03	0.00	0.03	0.00	0.10	0.20	0.30	0.57
	pre	0.07	0.03	0.13	0.10	0.13	0.33	0.47	0.60
l_img	raw	0.20	0.07	0.27	0.37	0.40	0.50	0.47	0.83
	pre	0.10	0.03	0.07	0.23	0.27	0.57	0.63	0.70
adv	raw	0.07	0.10	0.37	0.37	0.43	0.70	0.67	0.90
	pre	0.27	0.33	0.53	0.67	0.60	0.83	0.80	0.87
ko	raw	0.17	0.33	0.37	0.50	0.60	0.83	0.83	0.97
	pre	0.27	0.47	0.57	0.77	0.67	0.87	0.87	0.97
m_img+ko	raw	0.00	0.03	0.23	0.53	0.53	0.67	0.67	1.00
	pre	0.17	0.43	0.50	0.73	0.80	1.00	1.00	1.00
oz	raw	0.37	0.77	0.97	1.00	1.00	1.00	1.00	1.00
	pre	0.60	0.93	1.00	1.00	1.00	1.00	1.00	1.00

268 **Detecting distribution shift.** Rabanser et al. [2019] introduced a large benchmark for the de-  
 269 tection of distribution shifts. We repeat their experiments by considering the datasets MNIST  
 270 [LeCun et al., 2010] and CIFAR10 [Krizhevsky, 2009]. We consider sample sizes  $n, m \in$   
 271  $\{10, 20, 50, 100, 200, 500, 1000, 10000\}$ . Each shift is applied on a fraction  $\delta \in \{0.1, 0.5, 1.0\}$   
 272 of the second sample in different runs. We consider the following shifts: **Adversarial (adv)**: Turn  
 273 some images into adversarial examples via FGSM [Goodfellow et al., 2015]; **Knock-out (ko)**: Re-  
 274 move samples from class 0; **Gaussian noise (gn)**: Add gaussian noise to images with standard  
 275 deviation  $\sigma \in \{1, 10, 100\}$  (denoted  $s\_gn$ ,  $m\_gn$ , and  $l\_gn$ ); **Image (img)**: Natural shifts to images  
 276 through combinations of random rotations,  $(x, y)$ -axis-translation, as well as zoom-in with different  
 277 strength (denoted  $s\_img$ ,  $m\_img$ , and  $l\_img$ ); **Image + knock-out ( $m\_img+ko$ )**: Fixed medium image  
 278 shift and a variable knock-out shift; **Only-zero + image ( $oz+m\_img$ )**: Only images from class 0 in  
 279 combination with a variable medium image shift. More details are given in [Rabanser et al., 2019].  
 280 In total, we run 33 different shift experiments on MNIST and CIFAR10 each and for each sample  
 281 size. Every setting is repeated for 5 times.

282 The methods of Rabanser et al. [2019] perform a dimensionality reduction by using the whole training  
 283 set (50.000 images for MNIST, 40.000 images for CIFAR10). The actual tests compare examples  
 284 from the validation set (10.000 images) to examples from the shifted test set (10.000 images). They  
 285 also consider a C2ST trained on the raw features, i.e. without seeing the whole training set.

286 We add four univariate AutoML witness tests: a) **AutoML (raw)** trains a regression model on the raw  
 287 data with MSE, which is our default, b) **AutoML (pre)** uses the same setting, but trains on the softmax  
 288 output of a pretrained classifier for MNIST/CIFAR10 respectively, which is the same representation  
 289 as BBSDs used, c) **AutoML (class)** trains a classifier and uses its predicted probabilities of class '1'  
 290 as witness function, d) **AutoML (bin)** uses the same as c) but only considers binary outputs.

291 Our findings are reported in Table 1. From Table 1a we see that AutoML (raw) achieves overall very  
 292 competitive performance in detecting the shifts, especially for large sample sizes. Moreover, we  
 293 see that AutoML (raw) and AutoML (class) achieve comparable performance which confirms our  
 294 findings of Remark 1. Thresholding the classification probabilities to binary outputs always harms  
 295 the performance, see AutoML (class) vs. AutoML (bin). We can also compare AutoML (bin) with  
 296 'classif', as reported by Rabanser et al. [2019]. While both use binary classifiers for the testing,  
 297 'classif' used a fixed architecture across all shifts. This illustrates the power of using AutoML, as  
 298 we find significantly better performance across all sample sizes. If instead of training on the raw  
 299 features we start from the ten dimensional pretrained features, i.e. AutoML (pre), the performance  
 300 is improved when the sample size is small. For large sample sizes, instead working with the raw  
 301 features gives higher power.

302 In Table 1b we report the test power for comparing AutoML (raw) with AutoML (pre) for the different  
303 shifts. Using the pretrained probabilities of the softmax output, it is extremely hard to detect Gaussian  
304 noise, while AutoML (raw) does a fairly good job here. This is consistent with the findings of  
305 Rabanser et al. [2019, Table 1(b)]. Apparently, the output probabilities of the pretrained models are  
306 quite invariant under small and medium noise on the inputs. For the other shifts, such as knock-outs,  
307 using the pretrained features improves performance, particularly at small sample sizes.

## 308 6 Discussion

309 **Bias-variance tradeoff.** Our results on the distribution shift benchmark indicate a bias-variance  
310 tradeoff when optimizing the witness in stage-I. Learning the witness function over a ten dimensional  
311 pretrained representation gives good test power for some shifts even for small sample sizes, however,  
312 at the cost of being almost unable to detect other shifts, such as local Gaussian noise. Thus, learning  
313 on pretrained features introduces a strong bias. On the other hand, learning directly on the raw  
314 features introduces little bias, even more so since we used AutoGluon’s TabularPredictor, which  
315 is not specifically designed for images. This has the effect that on small sample sizes the test power is  
316 reduced, but when large data is available, we observe good test power across almost all shifts. For  
317 practical applications this implies that using models with the right bias when learning hypothesis  
318 tests is just as important as in any other supervised learning setting.

319 **Stand on the shoulders of giants.** As we see from the Blob and Higgs experiments the conceptually  
320 simple witness two-sample test can outperform more sophisticated test statistics like the deep MMD.  
321 This is possible through both the use of cross-validation (kfda-witness) or a full AutoML pipeline. In  
322 the distribution shift benchmark, we saw much better performance even when comparing a binary  
323 classifier (AutoML (bin)) with a classifier having a prespecified architecture (classif). Furthermore,  
324 using an AutoML framework allows practitioners to stand on the shoulders of giants and removes the  
325 need for specialized expertise. Instead, they can directly control how much time and resources to  
326 spend on optimizing the witness, which can lead to improved significance and/or inference time.

327 **Which test to use?** Obviously, there is no general answer to this question, and we are not claiming  
328 that our AutoML two-sample test should always be used. In special settings, a simple parametric test  
329 would perform much better than our AutoML witness test. Similarly, using MMD with a kernel can  
330 be the right choice in some settings. Nevertheless, a few points should be considered. For example,  
331 we demonstrate that a test using binary outputs of a classifier underperforms a test using the predicted  
332 probabilities of the same classifier. Therefore, we do recommend choosing the latter instead of the  
333 former. Furthermore, when using data splitting we should ensure that in the first stage we are actually  
334 optimizing the test power or a directly related proxy loss. To this end, it is important to use techniques  
335 that ensure good predictive performance and prevent overfitting. This brings us to the last point: We  
336 should also consider the resources available, both computational and human, that are relevant when  
337 implementing the test. That is, a testing framework should be easy to apply by a large group of users  
338 and should be adaptable to the computational resources the user is willed to spend on the test. The  
339 AutoML witness test can tick off all boxes. It learns a continuous witness function to optimize test  
340 power, leverages well-engineered toolboxes to maximize predictive performance, and requires little  
341 engineering expertise to apply and gives easy control over the computational resources used to learn  
342 the test, by setting a time limit and providing the available hardware.

## 343 7 Conclusion

344 We showed that optimizing a squared loss or cross-entropy loss leads to a witness function that  
345 maximizes test power, when using the mean discrepancy of the witness as a test statistic. This allows  
346 us to harness the advances in Automated Machine Learning, where regression and classification are  
347 the standard tasks, for two-sample testing. Although less studied, the use of a well-engineered toolbox  
348 to maximize the predictive performance of the learned function is as just as important for hypothesis  
349 testing as it is for supervised learning tasks. The result is a testing pipeline that is theoretically  
350 justified, leads to competitive performance, and is simple to apply in various settings. Our work thus  
351 constitutes a step towards fully automated statistical analysis of complex data [Steinruecken et al.,  
352 2019].

353 **References**

- 354 P. Baldi, P. Sadowski, and D. Whiteson. Searching for exotic particles in high-energy physics with  
355 deep learning. *Nature communications*, 5(1):1–9, 2014.
- 356 H. Cai, B. Goggin, and Q. Jiang. Two-sample test based on classification probability. *Statistical*  
357 *Analysis and Data Mining: The ASA Data Science Journal*, 13(1):5–13, 2020.
- 358 X. Cheng and A. Cloninger. Classification logit two-sample testing by neural networks.  
359 *arXiv:1909.11298*, 2019.
- 360 K. Chwialkowski, A. Ramdas, D. Sejdinovic, and A. Gretton. Fast two-sample testing with analytic  
361 representations of probability measures. In *NeurIPS*, 2015.
- 362 T. G. Dietterich. Ensemble methods in machine learning. In *International workshop on multiple*  
363 *classifier systems*, pages 1–15. Springer, 2000.
- 364 R. O. Duda, P. E. Hart, and D. G. Stork. *Pattern classification, 2nd Edition*. Wiley, 2001.
- 365 N. Erickson, J. Mueller, A. Shirkov, H. Zhang, P. Larroy, M. Li, and A. Smola. Autogluon-tabular:  
366 Robust and accurate automl for structured data. *arXiv:2003.06505*, 2020.
- 367 M. Feurer, A. Klein, J. Eggenberger, Katharina Springenberg, M. Blum, and F. Hutter. Efficient and  
368 robust automated machine learning. In *NeurIPS*, 2015.
- 369 J. H. Friedman. On multivariate goodness of fit and two sample testing. *Stanford Linear Accelerator*  
370 *Center-PUB-10325*, 2003.
- 371 P. Golland and B. Fischl. Permutation tests for classification: Towards statistical significance in image-  
372 based studies. In C. Taylor and J. A. Noble, editors, *Information Processing in Medical Imaging*,  
373 pages 330–341, Berlin, Heidelberg, 2003. Springer Berlin Heidelberg. ISBN 978-3-540-45087-0.
- 374 I. Goodfellow, J. Shlens, and C. Szegedy. Explaining and harnessing adversarial examples. In *ICLR*,  
375 2015.
- 376 A. Gretton, O. Bousquet, A. Smola, and B. Schölkopf. Measuring statistical dependence with  
377 hilbert-schmidt norms. In *International conference on algorithmic learning theory*, pages 63–77.  
378 Springer, 2005.
- 379 A. Gretton, K. M. Borgwardt, M. J. Rasch, B. Schölkopf, and A. Smola. A kernel two-sample test.  
380 *Journal of Machine Learning Research*, 13:723–773, 2012a.
- 381 A. Gretton, D. Sejdinovic, H. Strathmann, S. Balakrishnan, M. Pontil, K. Fukumizu, and B. K.  
382 Sriperumbudur. Optimal kernel choice for large-scale two-sample tests. In *NeurIPS*, 2012b.
- 383 X. He, K. Zhao, and X. Chu. Automl: A survey of the state-of-the-art. *Knowledge-Based Systems*,  
384 212:106622, 2021.
- 385 S. Hediger, L. Michel, and J. Näf. On the use of random forest for two-sample testing. *Computational*  
386 *Statistics and Data Analysis*, 170:107435, 2022.
- 387 F. Hutter, L. Kotthoff, and J. Vanschoren. *Automated machine learning: methods, systems, challenges*.  
388 Springer Nature, 2019.
- 389 W. Jitkrittum, Z. Szabó, K. P. Chwialkowski, and A. Gretton. Interpretable distribution features with  
390 maximum testing power. In *NeurIPS*, 2016.
- 391 I. Kim, A. Ramdas, A. Singh, and L. Wasserman. Classification accuracy as a proxy for two-sample  
392 testing. *The Annals of Statistics*, 49(1):411 – 434, 2021.
- 393 M. Kirchler, S. Khorasani, M. Kloft, and C. Lippert. Two-sample testing using deep learning. In  
394 *AISTATS*, 2020.
- 395 L. Koch, C. Schübach, A. Gretton, and P. Berens. Hidden in plain sight: Subgroup shifts escape ood  
396 detection. In *Medical Imaging in Deep Learning*, 2022.

- 397 A. Krizhevsky. Learning multiple layers of features from tiny images, 2009.
- 398 J. M. Kübler, W. Jitkrittum, B. Schölkopf, and K. Muandet. Learning kernel tests without data  
399 splitting. In *NeurIPS*, 2020.
- 400 J. M. Kübler, W. Jitkrittum, B. Schölkopf, and K. Muandet. A witness two-sample test. In *AISTATS*,  
401 2022.
- 402 Y. LeCun, C. Cortes, and C. Burges. Mnist handwritten digit database. *ATT Labs [Online]*. Available:  
403 <http://yann.lecun.com/exdb/mnist>, 2, 2010.
- 404 Z. Lipton, Y.-X. Wang, and A. Smola. Detecting and correcting for label shift with black box  
405 predictors. In *ICML*, 2018.
- 406 F. Liu, W. Xu, J. Lu, G. Zhang, A. Gretton, and D. J. Sutherland. Learning deep kernels for  
407 non-parametric two-sample tests. In *ICML*, 2020.
- 408 F. Liu, W. Xu, J. Lu, and D. J. Sutherland. Meta two-sample testing: Learning kernels for testing  
409 with limited data. In *NeurIPS*, 2021.
- 410 D. Lopez-Paz and M. Oquab. Revisiting classifier two-sample tests. In *ICLR*, 2017.
- 411 X. Mao, Q. Li, H. Xie, R. K. Lau, Z. Wang, and S. Smolley. On the effectiveness of least squares  
412 generative adversarial networks. *IEEE Transactions on Pattern Analysis and Machine Intelligence*,  
413 41(12):2947–2960, dec 2019.
- 414 S. Mika. *Kernel Fisher Discriminants*. Doctoral thesis, Technische Universität Berlin, Berlin, 2003.
- 415 S. Rabanser, S. Günnemann, and Z. Lipton. Failing loudly: An empirical study of methods for  
416 detecting dataset shift. In *NeurIPS*, 2019.
- 417 C. Steinruecken, E. Smith, D. Janz, J. Lloyd, and Z. Ghahramani. The automatic statistician. In  
418 *Automated Machine Learning*, pages 161–173. Springer, Cham, 2019.
- 419 Student. The probable error of a mean. *Biometrika*, pages 1–25, 1908.
- 420 D. J. Sutherland, H.-Y. Tung, H. Strathmann, S. De, A. Ramdas, A. Smola, and A. Gretton. Generative  
421 models and model criticism via optimized maximum mean discrepancy. In *ICLR*, 2017.
- 422 B. L. Welch. The generalization of ‘student’s’ problem when several different population variances  
423 are involved. *Biometrika*, 34(1-2):28–35, 1947.
- 424 S. Zhao, A. Sinha, Y. He, A. Perreault, J. Song, and S. Ermon. Comparing distributions by measuring  
425 differences that affect decision making. In *ICLR*, 2022.

## 426 Checklist

- 427 1. For all authors...
- 428 (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s  
429 contributions and scope? [Yes]
- 430 (b) Did you describe the limitations of your work? [Yes] Section 6 and Section 3.2
- 431 (c) Did you discuss any potential negative societal impacts of your work? [N/A]
- 432 (d) Have you read the ethics review guidelines and ensured that your paper conforms to  
433 them? [Yes]
- 434 2. If you are including theoretical results...
- 435 (a) Did you state the full set of assumptions of all theoretical results? [Yes]
- 436 (b) Did you include complete proofs of all theoretical results? [Yes] Partially in the  
437 Appendix.
- 438 3. If you ran experiments...

- 439 (a) Did you include the code, data, and instructions needed to reproduce the main ex-  
440 perimental results (either in the supplemental material or as a URL)? [Yes] Code is  
441 provided in the supplementary material and will be publicly available.
- 442 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they  
443 were chosen)? [Yes]
- 444 (c) Did you report error bars (e.g., with respect to the random seed after running experi-  
445 ments multiple times)? [Yes] In the tabular results, we omitted the error estimates for  
446 improved readability.
- 447 (d) Did you include the total amount of compute and the type of resources used (e.g., type  
448 of GPUs, internal cluster, or cloud provider)? [Yes] Section 5
- 449 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 450 (a) If your work uses existing assets, did you cite the creators? [Yes]
- 451 (b) Did you mention the license of the assets? [No]
- 452 (c) Did you include any new assets either in the supplemental material or as a URL? [No]
- 453 (d) Did you discuss whether and how consent was obtained from people whose data you're  
454 using/curating? [N/A]
- 455 (e) Did you discuss whether the data you are using/curating contains personally identifiable  
456 information or offensive content? [N/A]
- 457 5. If you used crowdsourcing or conducted research with human subjects...
- 458 (a) Did you include the full text of instructions given to participants and screenshots, if  
459 applicable? [N/A]
- 460 (b) Did you describe any potential participant risks, with links to Institutional Review  
461 Board (IRB) approvals, if applicable? [N/A]
- 462 (c) Did you include the estimated hourly wage paid to participants and the total amount  
463 spent on participant compensation? [N/A]