
Graph-based Discriminators: Sample Complexity and Expressiveness

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 A basic question in learning theory is to identify if two distributions are identical
2 when we have access only to examples sampled from the distributions. This
3 basic task is considered, for example, in the context of Generative Adversarial
4 Networks (GANs), where a discriminator is trained to distinguish between a real-
5 life distribution and a synthetic distribution. Classically, we use a hypothesis class
6 H and claim that the two distributions are distinct if for some $h \in H$ the expected
7 value on the two distributions is (significantly) different.

8 Our starting point is the following fundamental problem: "is having the hypothesis
9 dependent on more than a single random example beneficial". To address this
10 challenge we define k -ary based discriminators, which have a family of Boolean
11 k -ary functions \mathcal{G} . Each function $g \in \mathcal{G}$ naturally defines a hyper-graph, indicating
12 whether a given hyper-edge exists. A function $g \in \mathcal{G}$ distinguishes between two
13 distributions, if the expected value of g , on a k -tuple of i.i.d examples, on the two
14 distributions is (significantly) different.

15 We study the expressiveness of families of k -ary functions, compared to the classi-
16 cal hypothesis class H , which is $k = 1$. We show a separation in expressiveness
17 of $k + 1$ -ary versus k -ary functions. This demonstrate the great benefit of having
18 $k \geq 2$ as distinguishers.

19 For $k \geq 2$ we introduce a notion similar to the VC-dimension, and show that it
20 controls the sample complexity. We proceed and provide upper and lower bounds
21 as a function of our extended notion of VC-dimension.

22 1 Introduction

23 The task of discrimination consists of a *discriminator* that receives finite samples from two distribu-
24 tions, say p_1 and p_2 , and needs to certify whether the two distributions are distinct. Discrimination has
25 a central role within the framework of Generative Adversarial Networks [12], where a discriminator
26 trains a neural net to distinguish between samples from a real-life distribution and samples generated
27 synthetically by another neural network, called a *generator*.

28 A possible formal setup for discrimination identifies the discriminator with some distinguishing class
29 $\mathcal{D} = \{f : X \rightarrow \mathbb{R}\}$ of *distinguishing functions*. In turn, the discriminator wishes to find the best
30 $d \in \mathcal{D}$ that distinguishes between the two distributions. Formally, she wishes to find $d \in \mathcal{D}$ such that¹

$$31 \left| \mathbb{E}_{x \sim p_1} [d(x)] - \mathbb{E}_{x \sim p_2} [d(x)] \right| > \sup_{d^* \in \mathcal{D}} \left| \mathbb{E}_{x \sim p_1} [d^*(x)] - \mathbb{E}_{x \sim p_2} [d^*(x)] \right| - \epsilon. \quad (1)$$

¹Note that with such d at hand, with an order of $O(1/\epsilon^2)$ examples one can verify if any discriminator in the class certifies that the two distributions are distinct.

32 For examples, in GANs, the class of distinguishing functions we will consider could be the class of
 33 neural networks trained by the discriminator.

34 The first term in the RHS of eq. (1) is often referred to as the *Integral Probability Metric* (IPM distance)
 35 w.r.t a class \mathcal{D} [17], denoted $\text{IPM}_{\mathcal{D}}$. As such, we can think of the discriminator as computing the
 36 $\text{IPM}_{\mathcal{D}}$ distance.

37 Whether two, given, distributions can be distinguished by the discriminator becomes, in the IPM
 38 setup, a property of the distinguishing class. Also, the number of examples needed to be observed
 39 will depend on the class in question. Thus, if we take a large expressive class of distinguishers, the
 40 discriminator can potentially distinguish between any two distributions that are far in total variation.
 41 In that extreme, though, the class of distinguishers would need to be very large and in turn, the
 42 number of samples needed to be observed scales accordingly. One could also choose a “small” class,
 43 but at a cost of smaller distinguishing power that yields smaller IPM distance.

44 For example, consider two distributions over $[n]$ to be distinguished. We could choose as a distin-
 45 guishing class the class of *all* possible subsets over n . This distinguishing class give rise to the total
 46 variation distance, but the sample complexity turns out to be $O(n)$. Alternatively we can consider the
 47 class of *singletons*: This class will induce a simple IPM distance, with graceful sample complexity,
 48 however in worst case the IPM distance can be as small as $O(1/n)$ even though the total variation
 49 distance is large.

50 Thus, IPM framework initiates a study of generalization complexity where we wish to understand
 51 what is the expressive power of each class and what is its sample complexity.

52 For this special case that \mathcal{D} consists of Boolean functions, the problem turns out to be closely related
 53 to the classical statistical learning setting and prediction [22]. The sample complexity (i.e., number of
 54 samples needed to be observed by the discriminator) is governed by a combinatorial measure termed
 55 *VC dimension*. Specifically, for the discriminator to be able to find a d as in eq. (1), she needs to
 56 observe order of $\Theta(\frac{\rho}{\epsilon^2})$ examples, where ρ is the VC dimension of the class \mathcal{D} [5, 22].

57 In this work we consider a natural extension of this framework to more sophisticated discriminators:
 58 For example, consider a discriminator that observes pairs of points from the distribution and checks
 59 for collisions – such a distinguisher cannot a priori be modeled as a test of Boolean functions, as the
 60 tester measures a relation between two points and not a property of a single point. The collision test
 61 has indeed been used, in the context of synthetic data generation, to evaluate the *diversity* of the
 62 synthetic distribution [2].

63 More generally, suppose we have a class of 2-ary Boolean functions: $\mathcal{G} = \{g : g(x_1, x_2) \rightarrow \{0, 1\}\}$
 64 and the discriminator wishes to (approximately) compute

$$\sup_{g \in \mathcal{G}} \left| \mathbb{E}_{(x_1, x_2) \sim p_1^2} [g(x_1, x_2)] - \mathbb{E}_{(x_1, x_2) \sim p_2^2} [g(x_1, x_2)] \right|. \quad (2)$$

65 Here p^2 denotes the product distribution over p . More generally, we may consider k -ary mappings,
 66 but for the sake of clarity, we will restrict our attention in this introduction to $k = 2$.

67 Such 2-ary Boolean mapping can be considered as graphs where $g(x_1, x_2) = 1$ symbolizes that there
 68 exists an edge between x_1 and x_2 and similarly $g(x_1, x_2) = 0$ denotes that there is no such edge.
 69 The collision test, for example, is modelled by a graph that contains only self-loops. We thus call
 70 such multi-ary statistical tests *graph-based distinguishers*.

71 Two natural question then arise

- 72 1. Do graph-based discriminators have any added distinguishing power over classical discrimi-
 73 nators?
- 74 2. What is the sample complexity of graph-based discriminators?

75 With respect to the first question we give an affirmative answer and we show a separation between the
 76 distinguishing power of graph-based discriminators and classical discriminators. As to the second
 77 question, we introduce a new combinatorial measure (termed *graph VC dimension*) that governs the
 78 sample complexity of graph-based discriminators – analogously to the VC characterization of the
 79 sample complexity of classical discriminators. We next elaborate on each of these two results.

80 As to the distinguishing power of graph-based discriminators, we give an affirmative answer in the
81 following sense: We show that there exists a single graph g such that, for any distinguishing class \mathcal{D}
82 with bounded VC dimension, and ϵ , there are two distributions p_1 and p_2 that are \mathcal{D} -indistinguishable
83 but g certifies that p_1 and p_2 are distinct. Namely, the quantity in eq. (2) is at least $1/4$ for $\mathcal{G} = \{g\}$.

84 This result may be surprising. It is indeed known that for any two distributions that are ϵ -far in total
85 variation, there exists a boolean mapping d that distinguishes between the two distributions. In that
86 sense, distinguishing classes are known to be universal. Thus, asymptotically, with enough samples
87 any two distribution can be ultimately distinguished via a standard distinguishing function.

88 Nevertheless, our result shows that, given finite data, the restriction to classes with finite capacity
89 is limiting, and there could be graph-based distinguishing functions whose distinguishing power is
90 not comparable to *any* class with finite capacity. We stress that the same graph competes with *all*
91 finite-capacity classes, irrespective of their VC dimension.

92 With respect to the second question, we introduce a new VC-like notion termed *graph VC dimension*
93 that extends naturally to graphs (and hypergraphs). On a high level, we show that for a class of graph-
94 based distinguishers with graph VC dimension ρ , $O(\rho)$ examples are sufficient for discrimination
95 and that $\Omega(\sqrt{\rho})$ examples are necessary. This leaves a gap of factor $\sqrt{\rho}$ which we leave as an open
96 question.

97 The notion we introduce is strictly weaker than the standard VC-dimension of families of multi-ary
98 functions, and the proofs we provide do not follow directly from classical results on learnability of
99 finite VC classes [22, 5]. In more details, a graph-based distinguishing class \mathcal{G} is a family of Boolean
100 functions over the product space of vertices \mathcal{V} : $\mathcal{G} \subseteq \{0, 1\}^{\mathcal{V}^2}$. As such it is equipped with a VC
101 dimension, the largest set of pairs of vertices that is shattered by \mathcal{G} .

102 It is not hard to show that finite VC is sufficient to achieve finite sample complexity bounds over 2-ary
103 functions [9]. It turns out, though, that it is not a necessary condition: For example, one can show
104 that the class of k -regular graphs has finite graph VC dimension but infinite VC dimension. Thus,
105 even though they are not learnable in the standard PAC setting, they have finite sample complexity
106 within the framework of discrimination.

107 The reason for this gap, between learnability and discriminability, is that learning requires uniform
108 convergence with respect to any possible distribution over pairs, while discrimination requires
109 uniform convergence only with respect to product distributions – formally then, it is a weaker task,
110 and, potentially, can be performed even for classes with infinite VC dimension.

111 1.1 Related Work

112 The task of discrimination has been considered as early as the work of Vapnik and Chervonenkis in
113 [22]. In fact, even though Vapnik and Chervonenkis original work is often referred in the context of
114 prediction, the original work considered the question of when the empirical frequency of Boolean
115 functions converges uniformly to the true probability over a class of functions. In that sense, this
116 work can be considered as a natural extension to k -ary functions and generalization of the notion of
117 VC dimension.

118 The work of [9, 8] studies also a generalization of VC theory to multi-ary functions, in the context
119 of ranking tasks and U-statistics. They study the standard notion of VC dimension. Specifically
120 they consider the function class as Boolean functions over multi-tuples and the VC dimension is
121 defined by the largest set of multi-tuples that can be shattered. Their work provides several interesting
122 fast-rate convergence guarantees. As discussed in the introduction, our notion of capacity is weaker,
123 and in general the results are incomparable.

124 **GANs** A more recent interest in discrimination tasks is motivated by the framework of GANs,
125 where a neural network is trained to distinguish between two sets of data – one is real and the other is
126 generated by another neural network called *generator*. Multi-ary tests have been proposed to assess
127 the quality of GANs networks. [2] suggests birthday paradox to evaluate *diversity* in GANs. [19]
128 uses Binning to assess the solution proposed by GANs.

129 Closer to this work [15] suggests the use of a discriminator that observes samples from the m -th
130 product distribution. Motivated by the problem of *mode collapse* they suggest a theoretical framework
131 in which they study the algorithmic benefits of such discriminators and observe that they can

132 significantly reduce mode collapse. In contrast, our work is less concerned with the problem of mode
 133 collapse directly and we ask in general if we can boost the distinguishing power of discriminators via
 134 multi-ary discrimination. Moreover, we provide several novel sample complexity bounds.

135 **Property Testing** A related problem to ours is that of testing closeness of distributions [3, 11].
 136 Traditionally, testing closeness of distribution is concerned with evaluating if two discrete distributions
 137 are close vs. far/identical in *total variation*. [11], motivated by graph expansion test, propose a
 138 collision test to verify if a certain distribution is close to uniform. Interestingly, a collision test is a
 139 graph-based discriminator which turns out to be optimal for the setting[18]. Our sample-complexity
 140 lower bounds are derived from these results. Specifically we reduce discrimination to testing
 141 uniformity [18]. Other lower bounds in the literature can be similarly used to achieve alternative
 142 (yet incomparable bounds) (e.g. [7] provides a $\Omega(n^{2/3}/\epsilon^{3/4})$ lower bounds for testing whether two
 143 distributions are far or close).

144 In contrast with the aforementioned setup, here we do not measure distance between distributions
 145 in terms of total variation but in terms of an IPM distance induced by a class of distinguishers. The
 146 advantage of the IPM distance is that it sometimes can be estimated with limited amount of samples,
 147 while the total variation distance scales with the size of the support, which is often too large to allow
 148 estimation.

149 Several works do study the question of distinguishing between two distributions w.r.t a finite capacity
 150 class of tests, Specifically the work of [14] studies refutation algorithms that distinguish between
 151 noisy labels and labels that correlate with a bounded hypothesis class. [21] studies a closely related
 152 question in the context of realizable PAC learning. A graph-based discriminator can be directly turned
 153 to a refutation algorithm, and both works of [14, 21] show reductions from refutation to learning.
 154 In turn, the agnostic bounds of [14] can be harnessed to achieve lower bounds for graph-based
 155 discrimination. Unfortunately this approach leads to suboptimal lower bounds. It would be interesting
 156 to see if one can improve the guarantees for such reductions, and in turn exploit it for our setting.

157 2 Problem Setup

158 2.1 Basic Notations – Graphs and HyperGraphs

159 Recall that a k -hypergraph g consists of a set \mathcal{V}_g of *vertices* and a collection of non empty k -
 160 tuples over \mathcal{V} : $E_g \subseteq \mathcal{V}^k$, which are referred to as *hyperedges*. If $k = 2$ then g is called a graph.
 161 1-hypergraphs are simply identified as subsets over \mathcal{V} . We will normally use d to denote such
 162 1-hypergraphs and will refer to them as *distinguishers*. A distinguisher d can be identified with a
 163 Boolean function according to the rule: $d(x) = 1$ iff $x \in E_d$.

164 Similarly we can identify a k -hypergraph with a function $g : \mathcal{V}^k \rightarrow \{0, 1\}$. Namely, for any graph g
 165 we identify it with the Boolean function

$$g(v_1, \dots, v_k) = \begin{cases} 1 & (v_1, \dots, v_k) \in E_g \\ 0 & \text{else} \end{cases}$$

166 We will further simplify and assume that g is *undirected*, this means that for any permutation
 167 $\pi : [k] \rightarrow [k]$, we have that

$$g(v_{\pi(1)}, v_{\pi(2)}, \dots, v_{\pi(k)}) = g(v_1, \dots, v_k).$$

168 We will call undirected k -hypergraphs, k -distinguishers. A collection of k -distinguishers over a
 169 common set of vertices \mathcal{V} will be referred to as a *k -distinguishing class*. If $k = 1$ we will simply call
 170 such a collection a *distinguishing class*. For $k > 1$ we will normally denote such a collection with \mathcal{G}
 171 and for $k = 1$ we will often use the letter \mathcal{D} .

172 Next, given a distribution P over vertices and a k -hypergraph g let us denote as follows the frequency
 173 of an edge w.r.t P :

$$\mathbb{E}_P(g) = \mathbb{E}_{\mathbf{v}_{1:k} \sim P^k} [g(\mathbf{v}_{1:k})] = P^k [\{(\mathbf{v}_1, \dots, \mathbf{v}_k) : (\mathbf{v}_1, \dots, \mathbf{v}_k) \in E_g\}],$$

174 where we use the notation $\mathbf{v}_{1:t}$ in shorthand for the sequence $(\mathbf{v}_1, \dots, \mathbf{v}_t) \in \mathcal{V}^t$, and P^k denotes the
 175 product distribution of P k times.

176 Similarly, given a sample $S = \{v_i\}_{i=1}^m$ we denote the empirical frequency of an edge:

$$\mathbb{E}_S(g) = \frac{1}{m^k} \sum_{\mathbf{u}_{1:k} \in S^k} g(\mathbf{u}_{1:k}) = \frac{|\{(\mathbf{u}_1, \dots, \mathbf{u}_k) \in E_g : \forall i, \mathbf{u}_i \in S\}|}{m^k}$$

177 As a final set of notations: Given a k -hypergraph g a sequence $\mathbf{v}_{1:n}$ where $n < k$, we define a
 178 $k - n$ -distinguisher $g_{\mathbf{v}_{1:n}}$ as follows:

$$g_{\mathbf{v}_{1:n}}(\mathbf{u}_{1:k-n}) = g(\mathbf{v}_1, \dots, \mathbf{v}_n, \mathbf{u}_1, \dots, \mathbf{u}_{k-n}).$$

179 In turn, we define the following distinguishing classes: For every sequence $\mathbf{v}_{1:n}$, $n < k$, the
 180 distinguishing class $\mathcal{G}_{\mathbf{v}_{1:n}}$ is defined as follows:

$$\mathcal{G}_{\mathbf{v}_{1:n}} = \{g_{\mathbf{v}_{1:n}} : g \in \mathcal{G}\} \quad (3)$$

181 Finally, we point out that we will mainly be concerned with the case that $|\mathcal{V}| \leq \infty$ or $\mathcal{V} = \mathbb{N}$.
 182 However, all the results here can be easily extended to other domains as long as certain (natural)
 183 measurability assumptions are given to ensure that VC theory holds (see [22, 4]).

184 2.2 IPM distance

185 Given a class of distinguishers \mathcal{D} the induced IPM distance [17], denoted by $\text{IPM}_{\mathcal{D}}$, is a (pseudo)-
 186 metric between distributions over \mathcal{V} defined as follows

$$\text{IPM}_{\mathcal{D}}(p_1, p_2) = \sup_{d \in \mathcal{D}} |\mathbb{E}_{p_1}(d) - \mathbb{E}_{p_2}(d)| = \sup_{d \in \mathcal{D}} \left| \mathbb{E}_{v \sim p_1} [d(v)] - \mathbb{E}_{v \sim p_2} [d(v)] \right|.$$

187 The definition can naturally be extended to a general family of graphs, and we define:

$$\text{IPM}_{\mathcal{G}}(p_1, p_2) = \sup_{g \in \mathcal{G}} |\mathbb{E}_{p_1}(g) - \mathbb{E}_{p_2}(g)| = \sup_{g \in \mathcal{G}} \left| \mathbb{E}_{\mathbf{v}_{1:k} \sim p_1^k} [g(\mathbf{v}_{1:k})] - \mathbb{E}_{\mathbf{v}_{1:k} \sim p_2^k} [g(\mathbf{v}_{1:k})] \right|$$

188 Another metric we would care about is the *total variation metric*. Given two distributions p_1 and p_2
 189 the total variation distance is defined as:

$$\text{TV}(p_1, p_2) = \sup_E |p_1(E) - p_2(E)|$$

190 where $E \subseteq \mathcal{V}^{\{0,1\}}$ goes over all measurable events.

191 In contrast with an IPM distance, the total variation metric is indeed a metric and any two distributions
 192 $p_1 \neq p_2$ we have that $\text{TV}(p_1, p_2) > 0$. In fact, for every distinguishing class \mathcal{D} , $\text{IPM}_{\mathcal{D}} \preceq \text{TV}$.²

193 For finite classes of vertices \mathcal{V} , it is known that the total variation metric is given by

$$\text{TV}(p_1, p_2) = \frac{1}{2} \sum_{v \in \mathcal{V}} |p_1(v) - p_2(v)|.$$

194 Further, if we let $\mathcal{D} = P(\mathcal{V})$ the power set of \mathcal{V} we obtain

$$\text{IPM}_{P(\mathcal{V})}(p_1, p_2) = \text{TV}(p_1, p_2).$$

195 2.3 Discriminating Algorithms

196 **Definition 1.** Given a distinguishing class \mathcal{G} a \mathcal{G} -discriminating algorithm A with sample complexity
 197 $m(\epsilon, \delta)$ is an algorithm that receives as input two finite samples $S = (S_1, S_2)$ of vertices and outputs
 198 a hyper-graph $g_S^A \in \mathcal{G}$ such that:

²we use the notation $f_1 \preceq f_2$ to denote that for every x, y we have $f_1(x, y) \leq f_2(x, y)$.

199 If S_1, S_2 are drawn IID from some unknown distributions p_1, p_2 respectively and $|S_1|, |S_2| > m(\epsilon, \delta)$
 200 then w.p. $(1 - \delta)$ the algorithm's output satisfies:

$$|\mathbb{E}_{p_1}(g_S^A) - \mathbb{E}_{p_2}(g_S^A)| > \text{IPM}_{\mathcal{G}}(p_1, p_2) - \epsilon.$$

201 The sample complexity of a class \mathcal{G} is then given by the smallest possible sample complexity of a
 202 \mathcal{G} -discriminating algorithm A .

203 A class \mathcal{G} is said to be discriminable if it has finite sample complexity. Namely there exists a
 204 discriminating algorithm for \mathcal{G} with sample complexity $m(\epsilon, \delta) < \infty$.

205 **VC classes are discriminable** For the case $k = 1$, discrimination is closely related to PAC learning.
 206 It is easy to see that a proper learning algorithm for a class \mathcal{D} can be turned into a discriminator:
 207 Indeed, given access to samples from two distributions p_1 and p_2 we can provide a learner with
 208 labelled examples from a distribution p defined as follows: $p(y = 1) = p(y = -1) = \frac{1}{2}$ and
 209 $p(\cdot | y = 1) = p_1$, and $p(\cdot | y = -1) = p_2$. Given access to samples from p_1 and p_2 we can clearly
 210 generate IID samples from the distribution p . If, in turn, we provide a learner with samples from p
 211 and it outputs a hypothesis $d \in \mathcal{D}$ we have that (w.h.p):

$$\begin{aligned} |\mathbb{E}_{p_1}(d) - \mathbb{E}_{p_2}(d)| &= 2 \left| \frac{1}{2} \mathbb{E}_{(x,y) \sim p_1 \times \{1\}} [yd(x)] + \frac{1}{2} \mathbb{E}_{(x,y) \sim p_2 \times \{-1\}} [yd(x)] \right| \\ &= 2 \left| \mathbb{E}_{(x,y) \sim p} [yd(x)] \right| \\ &= 2(1 - 2p(d(x) \neq y)) \\ &\geq 2(1 - 2(\min_{d \in \mathcal{D}} p(d(x) \neq y) + \epsilon)) \\ &= \max_{d \in \mathcal{D}} (2 \left| \mathbb{E}_{(x,y) \sim p} yd(x) \right| - 4\epsilon) \\ &= \max_{d \in \mathcal{D}} |\mathbb{E}_{p_1}(d) - \mathbb{E}_{p_2}(d)| - 4\epsilon \\ &= \text{IPM}_{\mathcal{D}}(p_1, p_2) - 4\epsilon \end{aligned}$$

212 One can also see that a converse relation holds, if we restrict our attention to learning balanced labels
 213 (i.e., $p(y = 1) = p(y = -1)$). Namely, given labelled examples from some balanced distribution,
 214 the output of a discriminator is a predictor that competes with the class of predictors induced by \mathcal{D} .

215 Overall, the above calculation, together with Vapnik and Chervonenkis's classical result [22] shows
 216 that classes with finite VC dimension ρ are discriminable with sample complexity $O(\frac{\rho}{\epsilon^2})$.³ The
 217 necessity of finite VC dimension for agnostic PAC-learning was shown in [1]. Basically the same
 218 argument shows that given a class \mathcal{D} , $\tilde{\Omega}(\frac{\rho}{\epsilon^2})$ examples are necessary for discrimination. We next
 219 introduce a natural extension of VC dimension to hypergraphs, which will play a similar role.

220 2.4 VC Dimension of hypergraphs

221 We next define the notion of graph VC dimension for hypergraphs, as we will later see this notion
 222 indeed characterizes the sample complexity of discriminating classes, and in that sense it is a natural
 223 extension of the notion of VC dimension for hypotheses classes:

224 **Definition 2.** Given a family of k -hypergraphs, \mathcal{G} : The graph VC dimension of the class \mathcal{G} , denoted
 225 $\text{gVC}(\mathcal{G})$, is defined inductively as follows: For $k = 1$ $\text{gVC}(\mathcal{G})$ is the standard notion of VC dimension,
 226 i.e., $\text{gVC}(\mathcal{G}) = \text{VC}(\mathcal{G})$. For $k > 1$:

$$\text{gVC}(\mathcal{G}) = \max_{v \in \mathcal{V}} \{\text{gVC}(\mathcal{G}_v)\}$$

227 Roughly, the graph VC dimension of a hypergraph is given by the VC dimension of the induced
 228 classes of distinguishers via projections. Namely, we can think of the VC dimension of hypergraphs
 229 as the projected VC dimension when we fix all coordinates in an edge except for one.

³Recall that the VC dimension of a class \mathcal{D} is the largest set that can be shattered by \mathcal{D} where a set $S \subseteq \mathcal{V}$ is said to be shattered if \mathcal{D} restricted to S consists of $2^{|S|}$ possible Boolean functions.

230 **3 Main Results**

231 We next describe the main results of this work. The results are divided into two sections: For the first
 232 part we characterize the sample complexity of graph-based distinguishing class. The second part is
 233 concerned with the expressive/distinguishing power of graph-based discriminators. All proofs are
 234 provided in appendices B and C respectively.

235 **3.1 The sample complexity of graph-based distinguishing class**

236 We begin by providing upper bounds to the sample complexity for discrimination

237 **Theorem 1** (Sample Complexity – Upper Bound). *Let \mathcal{G} be a k -distinguishing class with $\text{gVC}(\mathcal{G}) =$
 238 ρ then \mathcal{G} has sample complexity $O(\frac{\rho k^2}{\epsilon^2} \log 1/\delta)$.*

239 Theorem 1 is a corollary of the following uniform convergence upper bound for graph-based distin-
 240 guishing classes.

241 **Theorem 2** (uniform convergence). *Let \mathcal{G} be a k -distinguishing class with $\text{gVC}(\mathcal{G}) = \rho$. Let
 242 $S = \{v_i\}_{i=1}^m$ be an IID sample of vertices drawn from some unknown distribution P . If $m =$
 243 $\Omega(\frac{\rho k^2}{\epsilon^2} \log 1/\delta)$ then with probability at least $(1 - \delta)$ (over the randomness of S):*

$$\sup_{g \in \mathcal{G}} |\mathbb{E}_S(g) - \mathbb{E}_P(g)| \leq \epsilon$$

244 The proof of theorem 2 is given in appendix B.1. We next provide a lower bound for the sample
 245 complexity of discriminating algorithms in terms of the graph VC dimension of the class

246 **Theorem 3** (Sample Complexity – Lower Bound). *Let \mathcal{G} be a k -distinguishing class with $\text{gVC}(\mathcal{G}) =$
 247 ρ . Any \mathcal{G} -discriminating algorithm with accuracy $\epsilon > 0$ that succeeds with probability $1 - \frac{2^{-k \log k}}{3}$,
 248 must observe at least $\Omega\left(\frac{\sqrt{\rho}}{2^{\tau k^3} \epsilon^2}\right)$ samples.*

249 We refer the reader to appendix B.2 for a proof of theorem 3. Our upper bounds and lower bounds
 250 leave a gap of order $O(\sqrt{\rho})$: As discussed in section 2.3, for the case $k = 1$ we can provide a tight
 251 $\theta(\frac{\rho}{\epsilon^2})$ bound through a reduction to agnostic pac learning and the appropriate lower bounds[1].

252 **3.2 The expressive power of graph-based distinguishing class**

253 So far we have characterized the discriminability of graph-based distinguishing classes. It is natural
 254 though to ask if graph-based distinguishing classes add any advantage over standard 1-distinguishing
 255 classes. In this section we provide several results that show that indeed graph provide extra expressive
 256 power over standard distinguishing classes.

257 We begin by providing a result over infinite graphs (proof is provided in appendix C.1)

258 **Theorem 4.** *Let $\mathcal{V} = \mathbb{N}$. There exists a distinguishing graph class \mathcal{G} , with sample complexity
 259 $m(\epsilon, \delta) = O(\frac{\log 1/\delta}{\epsilon^2})$ (in fact $|\mathcal{G}| = 1$) such that: for any 1-distinguishing class \mathcal{D} with finite VC
 260 dimension, and every $\epsilon > 0$ there are two distributions p_1, p_2 such that $\text{IPM}_{\mathcal{D}}(p_1, p_2) < \epsilon$ but
 261 $\text{IPM}_{\mathcal{G}}(p_1, p_2) > 1/2$*

262 Theorem 4 can be generalized to higher order distinguishing classes (see appendix C.2 for a proof):

263 **Theorem 5.** *Let $\mathcal{V} = \mathbb{N}$. There exists a k -distinguishing class \mathcal{G}_k , with sample complexity
 264 $m(\epsilon, \delta) = O(\frac{k^2 + \log 1/\delta}{\epsilon^2})$ such that: For any $k - 1$ -distinguishing class \mathcal{G}_{k-1} with bounded sample
 265 complexity, and every $\epsilon > 0$ there are two distributions p_1, p_2 such that $\text{IPM}_{\mathcal{G}_{k-1}}(p_1, p_2) < \epsilon$ and
 266 $\text{IPM}_{\mathcal{G}_k}(p_1, p_2) > 1/4$.*

267 **Finite Graphs** We next study the expressive power of distinguishing graphs over finite domains.

268 It is known that, over a finite domain $\mathcal{V} = \{1, \dots, n\}$, we can learn with a sample complexity of
 269 $O(\frac{n}{\epsilon^2} \log 1/\delta)$ any distinguishing class. In fact, we can learn the total variation metric (indeed the
 270 sample complexity of $\mathcal{P}(\mathcal{V})$ is bounded by $\log |\mathcal{P}(\mathcal{V})| = n$).

271 Therefore if we allow classes whose sample complexity scales linearly with n we cannot hope to
 272 show any advantage for distinguishing graphs. However, in most natural problems n is considered to

273 be very large (for example, over the Boolean cube n is exponential in the dimension). We thus, in
 274 general, would like to study classes that have better complexity in terms of n . In that sense, we can
 275 show that indeed distinguishing graphs yield extra expressive power.

276 In particular, we show that for classes with sublogarithmic sample complexity, we can construct
 277 graphs that are incomparable with a higher order distinguishing class.

278 **Theorem 6.** *Let $|\mathcal{V}| = n$. There exists a k -distinguishing class \mathcal{G}_k , with sample complexity $m(\epsilon, \delta) =$
 279 $O(\frac{k^2 + \log 1/\delta}{\epsilon^2})$ (in fact $|\mathcal{G}| = 1$) such that: For any $\epsilon > 0$ and any $k - 1$ distinguishing class \mathcal{G}_{k-1} if:*

$$\text{IPM}_{\mathcal{G}_{k-1}} \succ \epsilon \cdot \text{IPM}_{\mathcal{G}_k}$$

280 *then $\text{gVC}(\mathcal{G}_{k-1}) = \Omega(\frac{\epsilon^2}{k^2} \sqrt{\log n})$.*

281 The proof is given in appendix C.3. We can improve the bound in theorem 6 for the case $k = 1$ (see
 282 appendix C.4 for proof).

283 **Theorem 7.** *Let $|\mathcal{V}| = n$. There exists a 2-distinguishing class \mathcal{G} , with sample complexity $m(\epsilon, \delta) =$
 284 $O(\frac{\log 1/\delta}{\epsilon^2})$ (in fact $|\mathcal{G}| = 1$) such that: For any $\epsilon > 0$ and any distinguishing class \mathcal{D} if:*

$$\text{IPM}_{\mathcal{D}} \succ \epsilon \cdot \text{IPM}_{\mathcal{G}}$$

285 *then $\text{gVC}(\mathcal{D}) = \tilde{\Omega}(\epsilon^2 \log n)$.*

286 4 Discussion and open problems

287 In this work we developed a generalization of the standard framework of discrimination to graph-based
 288 distinguishers that discriminate between two distributions by considering multi-ary tests. Several
 289 open question arise from our results:

290 **Improving Sample Complexity Bounds** In terms of sample complexity, while we give a natural
 291 upper bound of $O(\rho k^2)$, the lower bound we provide are not tight neither in d nor in k and we provide
 292 a lower bound of $\Omega(\frac{\sqrt{\rho}}{2^{\text{poly}(k)}})$ This leave room for improvement both in terms of ρ and in terms of k .

293 **Improving Expressiveness Bounds** We also showed that, over finite domains, we can construct
 294 a graph that is incomparable with any class with VC dimension $\Omega(\epsilon^2 \log n)$. The best upper bound
 295 we can provide (the VC of a class that competes with any graph) is the naive $O(n)$ which is the VC
 296 dimension of the total variation metric.

297 Additionally, for the k -hypergraph case, our bounds deteriorate to $\Omega(\epsilon^2 \sqrt{\log n})$. The improvement in
 298 the graph case follows from using an argument in the spirit of Boosting [10] and Hardcore Lemma
 299 [13] to construct two indistinguishable probabilities with distinct support over a small domain. It
 300 would be interesting to extend these techniques in order to achieve similar bounds for the $k > 2$ case.

301 **Relation to GANs and Extension to Online Setting** Finally, a central motivation for learning the
 302 sample complexity of discriminators is in the context of GANs. It then raises interesting questions as
 303 to the *foolability* of graph-based distinguishers.

304 The work of [6] suggests a framework for studying sequential games between generators and
 305 discriminators (*GAM-Fooling*). In a nutshell, the GAM setting considers a sequential game between
 306 a generator G that outputs distributions and a discriminator D that has access to data from some
 307 distribution p^* (not known to G). At each round of the game, the generator proposes a distribution
 308 and the discriminator outputs a $d \in \mathcal{D}$ which distinguishes between the distribution of G and the true
 309 distribution p^* . The class \mathcal{D} is said to be GAM-Foolable if the generator outputs after finitely many
 310 rounds a distribution p that is \mathcal{D} -indistinguishable from p^*

311 [6] showed that a class \mathcal{D} is GAM-foolable if and only if it has finite Littlestone dimension. We then
 312 ask, similarly, which classes of graph-based distinguishers are GAM-Foolable? A characterization of
 313 such classes can potentially lead to a natural extension of the Littlestone notion and online prediction,
 314 to graph-based classes analogously to this work w.r.t VC dimension

References

- 315
- 316 [1] Martin Anthony and Peter L Bartlett. *Neural network learning: Theoretical foundations.*
317 cambridge university press, 2009.
- 318 [2] Sanjeev Arora and Yi Zhang. Do gans actually learn the distribution? an empirical study. *arXiv*
319 *preprint arXiv:1706.08224*, 2017.
- 320 [3] Tugkan Batu, Lance Fortnow, Ronitt Rubinfeld, Warren D Smith, and Patrick White. Testing
321 that distributions are close. In *Proceedings 41st Annual Symposium on Foundations of Computer*
322 *Science*, pages 259–269. IEEE, 2000.
- 323 [4] Shai Ben-David. 2 notes on classes with vapnik-chervonenkis dimension 1. *arXiv preprint*
324 *arXiv:1507.05307*, 2015.
- 325 [5] Anselm Blumer, Andrzej Ehrenfeucht, David Haussler, and Manfred K Warmuth. Learnability
326 and the vapnik-chervonenkis dimension. *Journal of the ACM (JACM)*, 36(4):929–965, 1989.
- 327 [6] Olivier Bousquet, Roi Livni, and Shay Moran. Passing tests without memorizing: Two models
328 for fooling discriminators. *arXiv preprint arXiv:1902.03468*, 2019.
- 329 [7] Siu-On Chan, Ilias Diakonikolas, Paul Valiant, and Gregory Valiant. Optimal algorithms for
330 testing closeness of discrete distributions. In *Proceedings of the twenty-fifth annual ACM-SIAM*
331 *symposium on Discrete algorithms*, pages 1193–1203. SIAM, 2014.
- 332 [8] Stéphan Cléménçon, Igor Colin, and Aurélien Bellet. Scaling-up empirical risk minimization:
333 optimization of incomplete u-statistics. *The Journal of Machine Learning Research*, 17(1):2682–
334 2717, 2016.
- 335 [9] Stéphan Cléménçon, Gábor Lugosi, Nicolas Vayatis, et al. Ranking and empirical minimization
336 of u-statistics. *The Annals of Statistics*, 36(2):844–874, 2008.
- 337 [10] Yoav Freund and Robert E Schapire. Game theory, on-line prediction and boosting. In *COLT*,
338 volume 96, pages 325–332. Citeseer, 1996.
- 339 [11] Oded Goldreich and Dana Ron. On testing expansion in bounded-degree graphs. In *Studies*
340 *in Complexity and Cryptography. Miscellanea on the Interplay between Randomness and*
341 *Computation*, pages 68–75. Springer, 2011.
- 342 [12] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil
343 Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial nets. In *Advances in neural*
344 *information processing systems*, pages 2672–2680, 2014.
- 345 [13] Russell Impagliazzo. Hard-core distributions for somewhat hard problems. In *Proceedings of*
346 *IEEE 36th Annual Foundations of Computer Science*, pages 538–545. IEEE, 1995.
- 347 [14] Pravesh K Kothari and Roi Livni. Agnostic learning by refuting. In *9th Innovations in*
348 *Theoretical Computer Science Conference (ITCS 2018)*. Schloss Dagstuhl-Leibniz-Zentrum
349 fuer Informatik, 2018.
- 350 [15] Zinan Lin, Ashish Khetan, Giulia Fanti, and Sewoong Oh. Pacgan: The power of two samples
351 in generative adversarial networks. In *Advances in Neural Information Processing Systems*,
352 pages 1498–1507, 2018.
- 353 [16] Richard J Lipton and Neal E Young. Simple strategies for large zero-sum games with appli-
354 cations to complexity theory. In *Proceedings of the twenty-sixth annual ACM symposium on*
355 *Theory of computing*, pages 734–740. ACM, 1994.
- 356 [17] Alfred Müller. Integral probability metrics and their generating classes of functions. *Advances*
357 *in Applied Probability*, 29(2):429–443, 1997.
- 358 [18] Liam Paninski. A coincidence-based test for uniformity given very sparsely sampled discrete
359 data. *IEEE Transactions on Information Theory*, 54(10):4750–4755, 2008.

- 360 [19] Eitan Richardson and Yair Weiss. On gans and gmms. In *Advances in Neural Information*
361 *Processing Systems*, pages 5847–5858, 2018.
- 362 [20] Shai Shalev-Shwartz and Shai Ben-David. *Understanding machine learning: From theory to*
363 *algorithms*. Cambridge university press, 2014.
- 364 [21] Salil P. Vadhan. On learning vs. refutation. In *Proceedings of the 30th Conference on Learning*
365 *Theory, COLT 2017, Amsterdam, The Netherlands, 7-10 July 2017*, pages 1835–1848, 2017.
- 366 [22] Vladimir N Vapnik and Aleksei Yakovlevich Chervonenkis. The uniform convergence of
367 frequencies of the appearance of events to their probabilities. In *Doklady Akademii Nauk*,
368 volume 181, pages 781–783. Russian Academy of Sciences, 1968.

369 A Prelimineries and Technical Background

370 A.1 Statistical Learning Theory

371 We begin with a brief overview of some classical results in Statistical Learning theory which
372 characterizes VC classes. Throughout we assume a domain \mathcal{X} and a *hypothesis class* which is a
373 family of Boolean functions over \mathcal{X} : $\mathcal{H} \subseteq \{0, 1\}^{\mathcal{X}}$.

374 **Theorem 8.** [Within proof of Thm. 6.11 in [20]] Let \mathcal{H} be a class with VC dimension ρ then

$$\mathbb{E}_{S \sim D^m} \left[\sup_{h \in \mathcal{H}} |\mathbb{E}_S(h) - \mathbb{E}_D(h)| \right] \leq \frac{4 + \sqrt{\rho \log(2em/\rho)}}{\sqrt{2m}}$$

375 Recall that a class \mathcal{H} has the *uniform convergence property*, if for some $m : (0, 1)^2 \rightarrow \mathbb{N}$ if P is some
376 unknown distribution and $S = \{x_i\}_{i=1}^m$ is a sample drawn IID from P such that $|S| > m(\epsilon, \delta)$ then
377 w.p. $(1 - \delta)$ (over the sample S):

$$\left| \frac{1}{m} \sum_{i=1}^m h(x_i) - \mathbb{E}_{x \sim P}[h(x)] \right| < \epsilon$$

378 The following, high probability analogue of theorem 8, is also an immediate corollary of Theorem
379 6.8 in [20]⁴:

380 **Corollary 1.** [Within Thm 6.8 [20]] Let \mathcal{D} be a class with VC dimension ρ . There exists a constant
381 $C > 0$, such that:

382 Let p be a distribution with finite support over \mathcal{V} . Let S be an IID sequence of m elements drawn
383 from p , and denote by p_S the empirical distribution over S . If $m \geq C \frac{\rho + \log 1/\delta}{\epsilon^2}$ then w.p. $(1 - \delta)$
384 (over the random choice of S) we have that

$$\text{IPM}_{\mathcal{D}}(p, p_S) = \sup_{d \in \mathcal{D}} |\mathbb{E}_p(d) - \mathbb{E}_{p_S}(d)| < \epsilon$$

385 A.2 Closeness Testing for Discrete Distribution

386 The problem of testing the closeness of two discrete distributions can be phrased as follows: Given
387 samples from two distributions p_1 and p_2 the tester needs to distinguish between the case $p_1 = p_2$
388 and the case that $\|p_1 - p_2\|_1 \geq \epsilon$. We will rely on the following result due to [18] (see also [7] for
389 discussion).

390 **Theorem 9.** Given $\epsilon > 0$ and access to samples from distributions p_1 and p_2 over $[n]$ any algorithm
391 that returns with probability $2/3$ *EQUIVALENT* if $p_1 = p_2$ and returns *DISTINCT* if
392 $\|p_1 - p_2\|_1 > \epsilon$ must observe at least $\Omega(\sqrt{n}/\epsilon^2)$ samples.

393 We note that [7] gives a slightly better lower bound, of $\Omega(\max(n^{3/4}/\epsilon^{4/3}, \sqrt{n}/\epsilon^2))$. However, our
394 proofs exploit other processes that exploit concentration inequalities, and it will be simpler to focus
395 on rates of order $O(1/\epsilon^2)$.

⁴Note that Theorem 6.8 is stated for $0 - 1$ loss, however considering a distribution with constant label $y = 0$
we can reduce the result for the loss $\ell(h, x) = h(x)$

396 B Sample Complexity –Proofs

397 B.1 Proof of theorem 2

398 **Theorem 2** (uniform convergence). *Let \mathcal{G} be a k -distinguishing class with $\text{gVC}(\mathcal{G}) = \rho$. Let*
 399 *$S = \{v_i\}_{i=1}^m$ be an IID sample of vertices drawn from some unknown distribution P . If $m =$*
 400 *$\Omega(\frac{\rho k^2}{\epsilon^2} \log 1/\delta)$ then with probability at least $(1 - \delta)$ (over the randomness of S):*

$$\sup_{g \in \mathcal{G}} |\mathbb{E}_S(g) - \mathbb{E}_P(g)| \leq \epsilon$$

401 Fix a k -distinguishing class \mathcal{G} with graph VC dimension ρ . As in the standard proof of uniform
 402 convergence for VC classes, we first prove the statement in expectation and then apply Mcdiarmid's
 403 inequality to prove the result w.h.p. Specifically, we will use the following Lemma (whose proof is
 404 given in appendix B.1.1):

405 **Lemma 1** (Uniform Convergence in Expectation). *Let \mathcal{G} be a k -distinguishing class with $\text{gVC}(\mathcal{G}) =$*
 406 *ρ . Let $S = \{v_i\}_{i=1}^m$ be an IID sample of vertices drawn from some unknown distribution P . Then,*

$$\mathbb{E}_{S \sim P^m} \left[\sup_{g \in \mathcal{G}} |\mathbb{E}_S(g) - \mathbb{E}_D(g)| \right] \leq \frac{k\sqrt{4 + \rho \log(2em/\rho)}}{\sqrt{2m}} + \frac{k(k-1)}{m}$$

407 We next proceed with the proof of theorem 2, assuming the correctness of lemma 1. Define

$$F(S) = \sup_{g \in \mathcal{G}} |\mathbb{E}_S(g) - \mathbb{E}_D(g)|,$$

408 Let $S = (v_1, \dots, v_m)$ be a sample and S' , some sequence that differ from S only in the i -th vertex
 409 then we will show that:

$$|F(S) - F(S')| \leq \frac{2k}{m} \tag{4}$$

410 Once we show eq. (4) holds, the result indeed follow from Mcdiarmid's inequality and lemma 1.
 411 Specifically if we assume that $m \geq \frac{8k^2(4 + \rho \log(2em/\rho))}{\epsilon^2} + \frac{2k^2 1/\delta}{\epsilon^2}$ then we obtain from lemma 1 that
 412 in expectation:

$$\mathbb{E}_{S \sim D^m} \sup_{g \in \mathcal{G}} |\mathbb{E}_S(g) - \mathbb{E}_D(g)| \leq \frac{\epsilon}{2}$$

413 Applying Mcdiarmid's we obtain that with probability at least $(1 - e^{-\frac{m\epsilon^2}{8k^2}})$, over the sample S :

$$F(S) - \mathbb{E}[F(S)] = \sup_{g \in \mathcal{G}} |\mathbb{E}_S(g) - \mathbb{E}_D(g)| - \mathbb{E}_{S \sim D^m} \sup_{g \in \mathcal{G}} |\mathbb{E}_S(g) - \mathbb{E}_D(g)| \leq \frac{\epsilon}{2}.$$

414 Noting that $m > \frac{8k^2 \log 1/\delta}{\epsilon^2}$, we obtain that with probability at least $(1 - \delta)$

$$F(S) = \sup_{g \in \mathcal{G}} |\mathbb{E}_S(g) - \mathbb{E}_D(g)| \leq \mathbb{E}_{S \sim D^m} \sup_{g \in \mathcal{G}} |\mathbb{E}_S(g) - \mathbb{E}_D(g)| + \frac{\epsilon}{2} \leq \epsilon$$

415 We are thus left with proving that eq. (4) holds.

416 For an index i and $m \geq i$, let us denote by $\pi_{i,m}$ all k -subsets of indices from $\{1, \dots, m\}$ that include
 417 i and we let $\pi_{-i,m}$ be all k -sequences that do not include i . Given a set S of size m let $S_{i,+}$ all the
 418 k -subsets of S that include v_i and let $S_{i,-}$ be all the k -subsets that do not include v_i . Next, denote

$$L_{S_{i,+}}(g) = \frac{1}{m^k} \sum_{(i_1, \dots, i_k) \in \pi_{i,m}} g(\mathbf{u}_{i_1}, \dots, \mathbf{u}_{i_k})$$

419 And similarly

$$L_{S_{i,-}}(g) = \frac{1}{m^k} \sum_{(i_1, \dots, i_k) \in \pi_{-i,m}} g(\mathbf{u}_{i_1}, \dots, \mathbf{u}_{i_k})$$

420 Then, let S and S' be two samples that differ on the i -th example. Specifically assume that $v_i \in S$
 421 and $v'_i \in S'$. Note that $S_{i,-} = S'_{i,-}$. Then:

$$\begin{aligned}
F(S) - F(S') &= \sup_{g \in \mathcal{G}} |\mathbb{E}_S(g) - \mathbb{E}_D(g)| - \sup_{g \in \mathcal{G}} |\mathbb{E}_{S'}(g) - \mathbb{E}_D(g)| \\
&= \sup_{g \in \mathcal{G}} |L_{S_{i,+}}(g) + L_{S_{i,-}}(g) - \mathbb{E}_D(g)| - \sup_{g \in \mathcal{G}} |L_{S'_{i,+}}(g) + L_{S'_{i,-}}(g) - \mathbb{E}_D(g)| \\
&\leq \sup_{g \in \mathcal{G}} |L_{S_{i,+}}(g) + L_{S_{i,-}}(g) - \mathbb{E}_D(g) - (L_{S'_{i,+}}(g) + L_{S'_{i,-}}(g) - \mathbb{E}_D(g))| \\
&= \sup_{g \in \mathcal{G}} |L_{S_{i,+}}(g) - L_{S'_{i,+}}(g)| \\
&= \sup_{g \in \mathcal{G}} |L_{S_{i,+}}(g)| + \sup_{g \in \mathcal{G}} |L_{S'_{i,+}}(g)| \\
&\leq \frac{|S_{i,+}|}{m^k} + \frac{|S'_{i,+}|}{m^k} \\
&= 2 \frac{m^k - (m-1)^k}{m^k} \\
&= 2 - 2\left(1 - \frac{1}{m}\right)^k \\
&\leq 2 \frac{k}{m}
\end{aligned}$$

422 We are thus left with proving lemma 1:

423 B.1.1 Proof of lemma 1

424 The proof of the statement follows by induction. The case $k = 1$ is the standard uniform convergence
425 property of VC classes, and it follows from theorem 8.

426 We next proceed to prove the statement for k , assuming it holds for $k - 1$. We begin by exploiting
427 trinagular inequality and together with adding/substracting terms:

$$\begin{aligned}
&\mathbb{E}_{S \sim D^m} \left[\sup_{g \in \mathcal{G}} |\mathbb{E}_S(g) - \mathbb{E}_D(g)| \right] \\
&= \mathbb{E}_{S \sim D^m} \left[\sup_{g \in \mathcal{G}} \left| \mathbb{E}_S(g) - \frac{1}{m^{k-1}} \sum_{\mathbf{v}_{1:k-1} \in S^{k-1}} \mathbb{E}_v g_{\mathbf{v}_{1:k-1}}(v) + \frac{1}{m^{k-1}} \sum_{\mathbf{v}_{1:k-1} \in S^{k-1}} \mathbb{E}_v g_{\mathbf{v}_{1:k-1}}(v) - \mathbb{E}_D(g) \right| \right] \\
&\leq \underbrace{\mathbb{E}_{S \sim D^m} \left[\sup_{g \in \mathcal{G}} \left| \mathbb{E}_S(g) - \frac{1}{m^{k-1}} \sum_{\mathbf{v}_{1:k-1} \in S^{k-1}} \mathbb{E}_v g_{\mathbf{v}_{1:k-1}}(v) \right| \right]}_* \\
&+ \\
&\underbrace{\mathbb{E}_{S \sim D^m} \left[\sup_{g \in \mathcal{G}} \left| \frac{1}{m^{k-1}} \sum_{\mathbf{v}_{1:k-1} \in S^{k-1}} \mathbb{E}_v g_{\mathbf{v}_{1:k-1}}(v) - \mathbb{E}_D(g) \right| \right]}_{**}
\end{aligned}$$

428 We next bound the two terms

Bounding *

$$\begin{aligned}
& \mathbb{E}_{S \sim D^m} \left[\sup_{g \in \mathcal{G}} \left| \frac{1}{m^{k-1}} \sum_{\mathbf{v}_{1:k-1} \in S^{k-1}} \frac{1}{m} \sum_{v \in S} g_{\mathbf{v}_{1:k-1}}(v) - \frac{1}{m^{k-1}} \sum_{\mathbf{v}_{1:k-1} \in S^{k-1}} \mathbb{E}_v g_{\mathbf{v}_{1:k-1}}(v) \right| \right] \\
& \leq \mathbb{E}_{S \sim D^m} \left[\frac{1}{m^{k-1}} \sum_{\mathbf{v}_{1:k-1} \in S^{k-1}} \sup_{d \in \mathcal{G}_{\mathbf{v}_{1:k-1}}} \left| \frac{1}{m} \sum_{v \in S} d(v) - \mathbb{E}_v d(v) \right| \right] \\
& = \mathbb{E}_{S \sim D^m} \left[\mathbb{E}_{\mathbf{v}_{1:k} \sim \mathcal{U}_{S^{k-1}}} \sup_{d \in \mathcal{G}_{\mathbf{v}_{1:k-1}}} \left| \frac{1}{m} \sum_{v \in S} d(v) - \mathbb{E}_v d(v) \right| \right]
\end{aligned}$$

429 where we denoted by $\mathcal{U}_{S^{k-1}}$ the uniform distribution over k -tuples from S . The expectation in the last
430 expression is thus taken w.r.t a process where we pick m elements according to d and then partition
431 them to $m - k + 1$ elements and to a sequence $\mathbf{v}_{1:k-1}$ of distinct element. This process is equivalent
432 to simply choosing $m - k + 1$ elements according to D , and then picking $k - 1$ new elements, again,
433 according to D . We thus continue and write:

$$\begin{aligned}
& = \mathbb{E}_{S \sim D^{m-k+1}} \mathbb{E}_{(\mathbf{v}_1, \dots, \mathbf{v}_{k-1}) \sim D^{k-1}} \left[\sup_{d \in \mathcal{G}_{\mathbf{v}_{1:k-1}}} \left| \frac{1}{m} \sum_{v \in S} d(v) + \frac{1}{m} \sum_{i=1}^{k-1} d(\mathbf{v}_i) - \mathbb{E}_v d(v) \right| \right] \\
& = \mathbb{E}_{S \sim D^{m-k+1}} \mathbb{E}_{(\mathbf{v}_1, \dots, \mathbf{v}_{k-1}) \sim D^{k-1}} \left[\sup_{d \in \mathcal{G}_{\mathbf{v}_{1:k-1}}} \left| \frac{1}{m} \sum_{v \in S} d(v) - \mathbb{E}_v d(v) + \frac{1}{m} \sum_{i=1}^{k-1} d(\mathbf{v}_i) \right| \right]
\end{aligned}$$

434 Note that the quantity $\frac{1}{m} \sum d(\mathbf{v}_i)$ is dependent on $\mathcal{G}_{\mathbf{v}_{1:k-1}}$, namely these are random sampled choices
435 that depend on our choice of distinguishing class. To bound their effect we next add and subtract
436 auxiliary random variables $\mathbf{u}_1, \dots, \mathbf{u}_{k-1}$ sampled IID according to D :

$$\begin{aligned}
& = \mathbb{E}_{S \sim D^{m-k+1}} \mathbb{E}_{(\mathbf{v}_1, \dots, \mathbf{v}_{k-1}) \sim D^{k-1}} \left[\sup_{d \in \mathcal{G}_{\mathbf{v}_{1:k-1}}} \left| \frac{1}{m} \sum_{v \in S} d(v) + \frac{1}{m} \mathbb{E}_{(\mathbf{u}_1, \dots, \mathbf{u}_{k-1}) \sim D^{k-1}} \sum d(\mathbf{u}_i) - \mathbb{E}_v d(v) \right. \right. \\
& \quad \left. \left. - \frac{1}{m} \mathbb{E}_{(\mathbf{u}_1, \dots, \mathbf{u}_k) \sim D^k} \sum d(\mathbf{u}_i) + \frac{1}{m} \sum_{i=1}^{k-1} d(\mathbf{v}_i) \right| \right] \\
& \leq \mathbb{E}_{S \sim D^{m-k+1}} \mathbb{E}_{(\mathbf{v}_1, \dots, \mathbf{v}_{k-1}) \sim D^{k-1}} \left[\sup_{d \in \mathcal{G}_{\mathbf{v}_{1:k-1}}} \left| \mathbb{E}_{(\mathbf{u}_1, \dots, \mathbf{u}_{k-1}) \sim D^{k-1}} \left[\frac{1}{m} \sum_{v \in S \cup \{\mathbf{u}_1, \dots, \mathbf{u}_{k-1}\}} d(v) \right] - \mathbb{E}_v [d(v)] \right| \right. \\
& \quad \left. + \left| \frac{1}{m} \mathbb{E}_{(\mathbf{u}_1, \dots, \mathbf{u}_k) \sim D^k} \sum d(\mathbf{u}_i) - \frac{1}{m} \sum_{i=1}^{k-1} d(\mathbf{v}_i) \right| \right] \\
& \leq \mathbb{E}_{(\mathbf{u}_1, \dots, \mathbf{u}_k) \sim D^k} \left[\mathbb{E}_{S \sim D^{m-k}} \mathbb{E}_{(\mathbf{v}_1, \dots, \mathbf{v}_{k-1}) \sim D^{k-1}} \left[\sup_{d \in \mathcal{G}_{\mathbf{v}_{1:k-1}}} \left| \frac{1}{m} \sum_{v \in S \cup \{\mathbf{u}_1, \dots, \mathbf{u}_k\}} d(v) - \mathbb{E}_v [d(v)] \right| \right] \right] + \frac{2k}{m}
\end{aligned}$$

437 Renaming $\mathbf{u}_1, \dots, \mathbf{u}_k$ and $\mathbf{v}_1, \dots, \mathbf{v}_k$ we can write:

$$\begin{aligned}
& \mathbb{E}_{(\mathbf{u}_1, \dots, \mathbf{u}_k) \sim D^k} \left[\mathbb{E}_{S \sim D^{m-k}} \mathbb{E}_{(\mathbf{v}_1, \dots, \mathbf{v}_{k-1}) \sim D^{k-1}} \left[\sup_{d \in \mathcal{G}_{\mathbf{v}_{1:k-1}}} \left| \frac{1}{m} \sum_{v \in S \cup \{\mathbf{u}_1, \dots, \mathbf{u}_k\}} d(v) - \mathbb{E}_v [d(v)] \right| \right] \right] + \frac{2k}{m} \\
& = \mathbb{E}_{(\mathbf{v}_1, \dots, \mathbf{v}_k) \sim D^k} \left[\mathbb{E}_{S \sim D^{m-k}} \mathbb{E}_{(\mathbf{u}_1, \dots, \mathbf{u}_{k-1}) \sim D^{k-1}} \left[\sup_{d \in \mathcal{G}_{\mathbf{u}_{1:k-1}}} \left| \frac{1}{m} \sum_{v \in S \cup \{\mathbf{v}_1, \dots, \mathbf{v}_k\}} d(v) - \mathbb{E}_v [d(v)] \right| \right] \right] + \frac{2k}{m} \\
& = \mathbb{E}_{(\mathbf{u}_1, \dots, \mathbf{u}_{k-1}) \sim D^{k-1}} \mathbb{E}_{S \sim D^m} \left[\sup_{d \in \mathcal{G}_{\mathbf{u}_{1:k-1}}} \left| \frac{1}{m} \sum_{v \in S} d(v) - \mathbb{E}_v [d(v)] \right| \right] + \frac{2k}{m}
\end{aligned}$$

438 Finally we apply theorem 8. Recalling that $\text{gVC}(\mathcal{D}_{\mathbf{u}_{1:k-1}}) = \rho$, and that the sequence S is drawn
 439 IID independent of the choice $\mathbf{u}_{1:k-1}$, we obtain for every fixed $(\mathbf{u}_1, \dots, \mathbf{u}_k)$

$$\mathbb{E}_{S \sim D^m} \left[\sup_{d \in \mathcal{D}_{\mathbf{u}_{1:k-1}}} \left| \frac{1}{m} \sum_{v \in S} d(v) - \mathbb{E}_v[d(v)] \right| \right] \leq \frac{4 + \sqrt{\rho \log 2em/\rho}}{\sqrt{2m}}$$

Bounding **

$$\begin{aligned} & \mathbb{E}_{S \sim D^m} \left[\sup_{g \in \mathcal{G}} \left| \frac{1}{m^{k-1}} \sum_{\mathbf{v}_{1:k-1} \in S^{k-1}} \mathbb{E}_v g_{\mathbf{v}_{1:k-1}}(v) - \mathbb{E}_{\mathbf{v}_{1:k-1}} \mathbb{E}_v g_{\mathbf{v}_{1:k-1}}(v) \right| \right] \\ & \leq \mathbb{E}_v \mathbb{E}_{S \sim D^m} \left[\sup_{g \in \mathcal{G}} \left| \frac{1}{m^{k-1}} \sum_{\mathbf{v}_{1:k-1} \in S^{k-1}} g_{\mathbf{v}_{1:k-1}}(v) - \mathbb{E}_{\mathbf{v}_{1:k-1}} g_{\mathbf{v}_{1:k-1}}(v) \right| \right] \\ & = \mathbb{E}_v \mathbb{E}_{S \sim D^m} \left[\sup_{g \in \mathcal{G}} \left| \frac{1}{m^{k-1}} \sum_{\mathbf{v}_{1:k-1} \in S^{k-1}} g_v(\mathbf{v}_1, \dots, \mathbf{v}_{k-1}) - \mathbb{E}_{\mathbf{v}_{1:k-1}} g_v(\mathbf{v}_1, \dots, \mathbf{v}_{k-1}) \right| \right] \\ & = \mathbb{E}_v \mathbb{E}_{S \sim D^m} \left[\sup_{g \in \mathcal{G}_v} \left| \frac{1}{m^{k-1}} \sum_{\mathbf{v}_{1:k-1} \in S^{k-1}} g(\mathbf{v}_1, \dots, \mathbf{v}_{k-1}) - \mathbb{E}_{\mathbf{v}_{1:k-1}} g(\mathbf{v}_1, \dots, \mathbf{v}_{k-1}) \right| \right] \end{aligned}$$

440 We now use the induction hypothesis: Note that \mathcal{G}_v is $(k-1)$ -distinguishing class with $\text{gVC}(\mathcal{G}_v) = \rho$
 441 for every choice of v . Thus, fixing v :

$$\begin{aligned} & \mathbb{E}_{S \sim D^m} \left[\sup_{g \in \mathcal{G}_v} \left| \frac{1}{m^{k-1}} \sum_{\mathbf{v}_{1:k-1} \in S^{k-1}} g(\mathbf{v}_1, \dots, \mathbf{v}_{k-1}) - \mathbb{E}_{\mathbf{v}_{1:k-1}} g(\mathbf{v}_1, \dots, \mathbf{v}_{k-1}) \right| \right] \\ & \leq \frac{(k-1) \left(4 + \sqrt{\rho \log(2em/\rho)} \right)}{\sqrt{2m}} + \frac{k(k-1)}{m} \end{aligned}$$

442 **Continuing the proof** With the aforementioned bound on the terms * and ** we now obtain

$$\begin{aligned} * + ** & \leq \frac{4 + \sqrt{\rho \log 2em/\rho}}{\sqrt{2m}} + \frac{2k}{m} + \frac{(k-1) \left(4 + \sqrt{\rho \log(2em/\rho)} \right)}{\sqrt{2m}} + \frac{k(k-1)}{m} \\ & = \frac{k \left(4 + \sqrt{\rho \log(2em/\rho)} \right)}{\sqrt{2m}} + \frac{(k+1)k}{m} \end{aligned}$$

443 B.2 Proof of theorem 3

444 **Theorem 3** (Sample Complexity – Lower Bound). *Let \mathcal{G} be a k -distinguishing class with $\text{gVC}(\mathcal{G}) =$
 445 ρ . Any \mathcal{G} -discriminating algorithm with accuracy $\epsilon > 0$ that succeeds with probability $1 - \frac{2^{-k \log k}}{3}$,
 446 must observe at least $\Omega\left(\frac{\sqrt{\rho}}{2^7 k^3 \epsilon^2}\right)$ samples.*

447 To prove theorem 3 we will in fact prove a stronger statement: We will show that it is not only hard
 448 to compute a $g \in \mathcal{G}$ as required, but in fact it is even hard to determine if such g exists vs. the case
 449 that $p_1 = p_2$.

450 Specifically let us call an algorithm A a testing algorithm for \mathcal{G} with sample complexity $m(\epsilon, \delta)$
 451 if A receives IID samples from two distributions p_1 and p_2 of size $m(\epsilon, \delta)$ and returns either
 452 *EQUIVALENT* or *DISTINCT* such that w.p. $(1 - \delta)$:

- 453 • If $p_1 = p_2$ the algorithm returns *EQUIVALENT* .

454 • If $\text{IPM}_{\mathcal{G}}(p_1, p_2) > \epsilon$ the algorithm returns *DISTINCT*

455 **Theorem 10.** Let \mathcal{G} be a k -distinguishing class with $\text{gVC}(\mathcal{G}) = \rho$. Any testing algorithm A with
 456 sample complexity $m(\epsilon, \delta)$ must observe $\Omega\left(\frac{\sqrt{\rho}}{2^{7k^3}\epsilon^2}\right)$ examples for any $\delta < \frac{2^{-k \log k}}{3}$.

457 Clearly, theorem 3 is a corollary of theorem 10. Indeed if A is a discriminating algorithm for \mathcal{G} with
 458 sample complexity $m(\epsilon, \delta)$ we can apply it over a sample of size $m(\epsilon/3, \delta)$ to receive (w.p. $1 - \delta$) a
 459 graph g s.t.

$$\text{IPM}_{\mathcal{G}}(p_1, p_2) \leq |\mathbb{E}_{p_1}(g) - \mathbb{E}_{p_2}(g)| + \frac{\epsilon}{3}.$$

460 With an additional sample of size $O\left(\frac{\log 1/\delta}{\epsilon^2}\right)$ we can estimate $|\mathbb{E}_{p_1}(g) - \mathbb{E}_{p_2}(g)|$ within accuracy $\epsilon/3$,
 461 and verify if $\text{IPM}_{\mathcal{G}}(p_1, p_2) < \epsilon$: The test will then output *EQUIVALENT* if $|\mathbb{E}_{p_1}(g) - \mathbb{E}_{p_2}(g)| <$
 462 $\frac{\epsilon}{3}$. It thus follows that, for sufficiently small δ $m(\epsilon, \delta) > \Omega\left(\frac{\sqrt{\rho}}{2^{7k^3}\epsilon^2}\right)$.

463 We proceed with the proof of theorem 10.

464 B.2.1 Proof of theorem 10

465 The proof is done by induction. For the induction, we will assume a more fine-grained lower bound.
 466 We will assume that there exists a constant C so that for every $n \leq k - 1$, if $m_n(\epsilon, \delta)$ is the sample
 467 complexity of a n -distinguishing class then:

$$m_n(\epsilon, \delta) \geq C \frac{\sqrt{\rho}}{(n+1)! 2^{\sum_{j=1}^n 6j^2} \cdot \epsilon^2} = \Omega\left(\frac{\sqrt{\rho}}{2^{7n^3}\epsilon^2}\right). \quad (5)$$

468 $C > 0$ will depend only on the constant for the lower bound for testing if two distributions are distinct
 469 or ϵ -far in total variation, as in theorem 9.

470 We start with the case $k = 1$.

471 $k = 1$ The case $k = 1$ follows directly from theorem 9. Let \mathcal{D} be a class with VC dimension ρ . by
 472 restricting our attention to probabilities supported on the shattered set of size ρ , we may assume that
 473 $|\mathcal{V}| = \rho$ and that $\mathcal{D} = P(\mathcal{V})$. Note then, that for the IPM distance we then have

$$\text{IPM}_{\mathcal{D}}(p_1, p_2) = \text{TV}(p_1, p_2).$$

474 theorem 9 immediately yields the result.

475 the induction step We now proceed with the proof assuming the statement holds for $k - 1$.

476 By assumption $\text{gVC}(\mathcal{G}) = \rho$. Fix $v \in \mathcal{V}$ such that $\text{gVC}(\mathcal{G}_v) = \rho$. For every $q \in (0, 1)$ and
 477 distribution p denote

$$p^q := q\delta_v + (1 - q)p. \quad (6)$$

478 We next state the core Lemma we will need for the proof:

479 **Lemma 2.** Let \mathcal{G} be a family of k -hypergraphs and p_1, p_2 two distributions. Assume that for some
 480 $v \in \mathcal{V}$ we have that:

$$\text{IPM}_{\mathcal{G}_v}(p_1, p_2) \geq \epsilon.$$

481 Let p_1^q and p_2^q be as in eq. (6) for our choice of $v \in \mathcal{V}$.

482 Then for some value $q \in \{0, \frac{1}{k}, \frac{2}{k}, \dots, 1\}$ we have that,

$$\text{IPM}_{\mathcal{G}}(p_1^q, p_2^q) \geq \frac{\epsilon}{2^{3k^2}}.$$

483 We deter the proof of lemma 2 to appendix B.2.2, and proceed with the proof of the induction step.
 484 Let us denote $\delta_k = 2^{-k \log k}$ and denote $c_k = 2^{-3k^2}$.

485 Let A be a testing algorithm for \mathcal{G} with sample complexity $m(\epsilon, \delta)$ as in theorem 10. We can now
 486 construct a testing algorithm for \mathcal{G}_v with sample complexity

$$m_{k-1}(\epsilon, \delta) = (k + 1) \cdot m\left(c_k \epsilon, \frac{\delta}{k}\right),$$

487 as follows: Run the testing algorithm A on pairs of distributions $(p_1, p_2), (p_1^{1/k}, p_2^{1/k}), \dots, (p_1^1, p_2^1)$,
 488 each on its own fixed sample of size $m(c_k \epsilon, \frac{\delta}{k})$. If the algorithm returns $DISTINCT$ for any of
 489 these tests, output $DISTINCT$, else output $EQUIVALENT$.

490 We now show that if $p_1 = p_2$ the algorithm outputs w.p. $(1 - \delta)$ $EQUIVALENT$: Indeed, since
 491 $p_1 = p_2$, we have that $p_1^q = p_2^q$ for all q : Applying union bound we have that w.p. $(1 - \delta)$ the
 492 algorithm indeed outputs $EQUIVALENT$.

493 On the other hand, if $\text{IPM}_{\mathcal{G}_v}(p_1, p_2) \geq \epsilon$ we have by lemma 2 that for one of the distributions
 494 (p_1^q, p_2^q) , $\text{IPM}_{\mathcal{G}}(p_1^q, p_2^q) > c_k \epsilon$, in particular the algorithm will output $DISTINCT$ with probability
 495 $(1 - \delta)$. Overall we constructed a testing algorithm for \mathcal{G}_v with sample complexity $(k + 1)m(c_k \epsilon, \frac{\delta}{k})$

496 Since $\delta_k < \frac{2^{-(k-1) \log(k-1)}}{k}$, it follows from the induction step

$$(k + 1)m(c_k \epsilon, \frac{\delta}{k}) = m_{k-1}(\epsilon, \delta) \geq C \frac{\sqrt{\rho}}{k! 2^{\sum_{n=1}^{k-1} 6n^2} \cdot \epsilon^2}$$

497 Reparametrizing we obtain

$$m(\epsilon, \frac{\delta}{k}) \geq C \frac{\sqrt{\rho}}{(k + 1)! 2^{\sum_{n=1}^k 6n^2} \cdot \epsilon^2}$$

498 B.2.2 Proof of lemma 2

499 Denote

$$\Delta_n^g(p_1, p_2) = \mathbb{E}_{\mathbf{u}_{1:n} \sim p_1^{k-n}} \underbrace{g(v, v, v, \dots, v, \mathbf{u}_1, \dots, \mathbf{u}_{k-n})}_n - \mathbb{E}_{\mathbf{u}_{1:n} \sim p_2^{k-n}} \underbrace{g(v, v, v, \dots, v, \mathbf{u}_1, \dots, \mathbf{u}_{k-n})}_n$$

500 One can show that

$$\begin{aligned} \text{IPM}_{\mathcal{G}}(p_1^q, p_2^q) &= \sup_{g \in \mathcal{G}} \left| \sum \binom{k}{n} q^n (1 - q)^{k-n} \Delta_n^g(p_1, p_2) \right| \\ &= \sup_{g \in \mathcal{G}} \left| (1 - q)^k \Delta_0^g(p_1, p_2) + kq(1 - q)^{n-1} \Delta_1^g(p_1, p_2) + \sum_{n=2}^k \binom{k}{n} q^n (1 - q)^{k-n} \Delta_n^g(p_1, p_2) \right| \\ &= \sup_{g \in \mathcal{G}} \left| \Delta_0^g(p_1, p_2) + kq(\Delta_1^g(p_1, p_2) - \Delta_0^g(p_1, p_2)) + q^2 p_g(q) \right| \end{aligned}$$

501 where $p_g(q)$ is some $k - 2$ degree polynomial in q whose coefficient depend on g and p_1 and p_2 . We
 502 next apply the following claim

503 **Claim 1.** Let $f(q) = a_0 + a_1 q + q^2 p(q)$ where $p(q)$ is some $k - 2$ degree polynomial. then for some
 504 value $q_0 \in \{0, \frac{1}{k}, \frac{2}{k}, \dots, 1\}$ we have that $|f(q_0)| \geq \frac{|a_1|}{2^{3k^2}}$

505 *Proof Sketch.* We provide a full proof for this claim in appendix D.1. In a nutshell, claim 1 follows
 506 from the equivalence between norms in finite dimensional spaces. Indeed, the mapping

$$(a_0, \dots, a_k) \rightarrow (p_a(1/k), p_a(2/k), \dots, p_a(1)),$$

507 where $p_a(x) = \sum a_i x^i$ is known to be a non-singular linear transformation induced by the appropri-
 508 ate Vandermonde matrix (specifically. $V_{i,j} = ((i - 1)/k)^{j-1}$). Letting λ_{min} be the smallest singular
 509 value of the matrix V , we know that $\|V\mathbf{a}\|_2 \geq \lambda_{min} \|\mathbf{a}\|_2$. where \mathbf{a} is the vector of coefficients of the
 510 polynomial p_a .

511 Finally, we exploit the relation in \mathbb{R}^{k+1} : $\|x\|_\infty \leq \|x\|_2 \leq \sqrt{k+1} \|x\|_\infty$. We can, thus, relate the
 512 max norm of the coefficient vector $\|a\|_\infty \geq |a_1|$ to the maximum value $\max_{i \in \{0, \dots, k\}} \sum a_j (i/k)^j =$
 513 $\|V\mathbf{a}\|_\infty$ to obtain

$$|a_1| \leq \|\mathbf{a}\|_2 \leq \lambda_{min}^{-1} \|V\mathbf{a}\|_2 \leq \frac{\sqrt{k+1}}{\lambda_{min}} \|V\mathbf{a}\|_\infty = \frac{\sqrt{k+1}}{\lambda_{min}} \max_{i \in \{0, \dots, k\}} \sum a_j (i/k)^j$$

514 It remains only to lower bound the singular values of V , this is done in the full proof in appendix D.1.
 515 □

516 With claim 1 in mind we prove the result as follows: First, suppose that for some $g \in \mathcal{G}$ we have
 517 that $|k(\Delta_0^g(p_1, p_2) - \Delta_1^g(p_1, p_2))| > \frac{\epsilon}{2}$. In this case, applying claim 1 with $a_0 = \Delta_0^g(p_1, p_2)$ and
 518 $a_1 = k(\Delta_0^g(p_1, p_2) - \Delta_1^g(p_1, p_2))$ and $p = p_g$, we obtain that there exists a value $q = j/k$ such that
 519 $\text{IPM}_{\mathcal{G}}(p_1^q, p_2^q) \geq \frac{\epsilon}{2^{3k^2}}$.

520 On the other hand, consider the case that $|k(\Delta_0^g(p_1, p_2) - \Delta_1^g(p_1, p_2))| < \frac{\epsilon}{2}$ for any $g \in \mathcal{G}$, by
 521 assumption we have that $|\Delta_1^g(p_1, p_2)| > \epsilon$, for some $g \in \mathcal{G}$. Hence $|\Delta_0^g(p_1, p_2)| > \epsilon$. By definition
 522 of Δ_0 we have that for $q = 0$ we obtain that: $\text{IPM}_{\mathcal{G}}(p_1^q, p_2^q) = |\mathbb{E}(p_1^q) - \mathbb{E}(p_2^q)| > \frac{\epsilon}{2}$.

523 C Expressivity – Proofs

524 C.1 Proof of theorem 4

525 **Theorem 4.** *Let $\mathcal{V} = \mathbb{N}$. There exists a distinguishing graph class \mathcal{G} , with sample complexity*
 526 $m(\epsilon, \delta) = O(\frac{\log 1/\delta}{\epsilon^2})$ (in fact $|\mathcal{G}| = 1$) such that: for any 1-distinguishing class \mathcal{D} with finite VC
 527 dimension, and every $\epsilon > 0$ there are two distributions p_1, p_2 such that $\text{IPM}_{\mathcal{D}}(p_1, p_2) < \epsilon$ but
 528 $\text{IPM}_{\mathcal{G}}(p_1, p_2) > 1/2$

529 As stated, the class \mathcal{G} will consist of a single graph g . The graph g is going to be a bipartite
 530 graph. We thus, divide the vertices into two infinite sets: \mathcal{V}_1 and \mathcal{V}_2 the elements of \mathcal{V}_1 will be
 531 indexed by \mathbb{N} i.e. $\mathcal{V}_1 = \{v_1, v_2, \dots\}$ and we index the elements of \mathcal{V}_2 with finite subsets of \mathbb{N}
 532 $\mathcal{V}_2 = \{v_A : A \subseteq \mathbb{N}, |A| < \infty\}$. Next we define g so that an edge passes between $v_i \in \mathcal{V}_1$ and
 533 $v_A \in \mathcal{V}_2$ iff $i \in A$.

534 Let \mathcal{D} be a distinguishing class with finite sample complexity, in particular $\text{gVC}(\mathcal{D}) < \infty$. Denote
 535 $\text{gVC}(\mathcal{D}) = \rho$. Let \mathcal{D}_1 be the restriction of \mathcal{D} to \mathcal{V}_1 : Note that $\text{gVC}(\mathcal{D}_1) \leq \rho$.

536 Next we make the following claim:

537 **Claim 2.** *There are two distributions, q_1 and q_2 , supported on \mathcal{V}_1 so that*

$$\text{IPM}_{\mathcal{D}_1}(p_1, p_2) < \epsilon.$$

538 *and yet q_1 and q_2 have disjoint support.*

539 *Proof.* To construct two such distributions, choose a set $S \subseteq \mathcal{V}_1$ of size m large enough (to be
 540 determined later). Then, randomly choose two samples S_1 and S_2 out of S (uniformly), each of size
 541 $O(\frac{\rho}{\epsilon^2})$. Then, by theorem 2 with some constant probability we have that $\text{IPM}_{\mathcal{D}}(p_{S_1}, p_S) < \epsilon/2$ and
 542 similarly $\text{IPM}_{\mathcal{D}}(p_S, p_{S_2}) < \epsilon/2$. Taken together we obtain that $\text{IPM}_{\mathcal{G}}(p_{S_1}, p_{S_2}) < \epsilon$.

543 Also, if S is sufficiently large (say, of order $O(\frac{\rho^2}{\epsilon^4})$), we would have that w.h.p $S_1 \cap S_2 = \emptyset$. Thus, let
 544 $q_1 = p_{S_1}$ and $q_2 = p_{S_2}$. □

545 With claim 2, we proceed with the proof. Let q_1 and q_2 be as in claim 2. Let A be the support of q_1 ,
 546 and define p_1 to be a distribution $p_1 = \frac{1}{2}\delta_A + \frac{1}{2}q_1$ and similarly we define $p_2 = \frac{1}{2}\delta_A + \frac{1}{2}q_2$. We then
 547 have

$$\begin{aligned} \text{IPM}_{\mathcal{D}}(p_1, p_2) &= \frac{1}{2}\text{IPM}_{\mathcal{D}}(q_1, q_2) \\ &= \frac{1}{2}\text{IPM}_{\mathcal{D}_1}(q_1, q_2) \\ &< \epsilon. \end{aligned}$$

548 On the other hand, note that for p_1 the probability to draw an edge from g is at least $1/2$ (indeed if
 549 $v_1 = v_A$ and $v_2 \neq v_A$ drawn from q_1 then $g(v_1, v_2) = 1$). On the other hand, the probability to draw
 550 an edge from p_2 is 0. It follows that

$$\text{IPM}_{\mathcal{G}}(p_1, p_2) > \frac{1}{2}.$$

551 **C.2 Proof of theorem 5**

552 **Theorem 5.** *Let $\mathcal{V} = \mathbb{N}$. There exists a k -distinguishing class \mathcal{G}_k , with sample complexity*
 553 *$m(\epsilon, \delta) = O(\frac{k^2 + \log 1/\delta}{\epsilon^2})$ such that: For any $k - 1$ -distinguishing class \mathcal{G}_{k-1} with bounded sample*
 554 *complexity, and every $\epsilon > 0$ there are two distributions p_1, p_2 such that $\text{IPM}_{\mathcal{G}_{k-1}}(p_1, p_2) < \epsilon$ and*
 555 *$\text{IPM}_{\mathcal{G}_k}(p_1, p_2) > 1/4$.*

556 The construction is similar to the case $k = 2$. We again divide the vertices into two infinite sets \mathcal{V}_1
 557 and \mathcal{V}_2 . Again, the elements of \mathcal{V}_1 will be indexed by \mathbb{N} , and the elements of \mathcal{V}_2 are indexed by finite
 558 subsets of \mathbb{N} . $\mathcal{V}_2 = \{v_A : A \subseteq \mathbb{N}, |A| < \infty\}$.

559 We define the hyper graph g_k to be a (undirected) graph that contains a hyperedge $(v_{i_1}, \dots, v_{i_{k-1}}, v_A)$
 560 whenever $\{i_1, \dots, i_{k-1}\} \subseteq A$.

561 Next, as before we construct two distributions with distinct support such that $\text{IPM}_{\mathcal{G}}(p_1, p_2) \leq \epsilon$.
 562 This is done similar to the proof of theorem 4. Specifically:

563 **Claim 3.** *Let \mathcal{G} be a $k - 1$ -distinguishing class defined on \mathcal{V}_1 . There are two distributions, q_1 and q_2 ,*
 564 *supported on \mathcal{V}_1 so that*

$$\text{IPM}_{\mathcal{G}}(p_1, p_2) < \epsilon.$$

565 *and yet q_1 and q_2 have disjoint support.*

566 The proof is a repetition of the proof of claim 2, where we draw S_1 and S_2 to be order of $O(\frac{k^2 \rho}{\epsilon^2})$,
 567 and again invoke theorem 2.

568 As before, then, given a class \mathcal{G} of $k - 1$ -hypergraphs we take two distributions q_1 and q_2 as in
 569 claim 3 and if A is the support of q_1 , we take $p_1 = \frac{1}{k} \delta_{v_A} + (1 - \frac{1}{k}) q_1$ and let $p_2 = \frac{1}{k} \delta_{v_A} + (1 - \frac{1}{k}) q_2$.
 570 Then, we can show that $\text{IPM}_{\mathcal{G}}(p_1, p_2) \leq \epsilon$. On the other hand, the probability to draw an edge from
 571 g_k is $k \cdot \frac{1}{k} (1 - \frac{1}{k})^{k-1} \geq e^{-1}$ according to p_1 , but the probability to draw an edge from p_2 is 0.

572 **C.3 Proof of theorem 6**

573 **Theorem 6.** *Let $|\mathcal{V}| = n$. There exists a k -distinguishing class \mathcal{G}_k , with sample complexity $m(\epsilon, \delta) =$
 574 *$O(\frac{k^2 + \log 1/\delta}{\epsilon^2})$ (in fact $|\mathcal{G}| = 1$) such that: For any $\epsilon > 0$ and any $k - 1$ distinguishing class \mathcal{G}_{k-1} if:**

$$\text{IPM}_{\mathcal{G}_{k-1}} \succ \epsilon \cdot \text{IPM}_{\mathcal{G}_k}$$

575 *then $\text{gVC}(\mathcal{G}_{k-1}) = \Omega(\frac{\epsilon^2}{k^2} \sqrt{\log n})$.*

576 The proof is similar to the proof of theorem 5. For simplicity, let us assume that $|\mathcal{V}| = n + \log n$.
 577 This will not change the results up to constants.

578 Given $n + \log n$ vertices we partition them into two sets \mathcal{V}_1 , of size $\log n$ and \mathcal{V}_2 . We index
 579 the elements of \mathcal{V}_1 as $\{v_1, \dots, v_{\log n}\}$ and we index the elements of \mathcal{V}_2 with subsets of $[\log n]$.
 580 We then consider a graph g that contains only hyper-edges of the form $(v_{i_1}, \dots, v_{i_{k-1}}, v_A)$ iff
 581 $\{i_1, \dots, i_{k-1}\} \in A$.

582 Next, let \mathcal{G}_{k-1} be a distinguishing class with $\text{gVC}(\mathcal{G}_{k-1}) = \rho$, and let $m(\epsilon, \delta) = O(\frac{\rho k^2}{\epsilon^2})$ be an
 583 upper bound on the sample complexity of classes of graph VC dimension ρ .

584 We claim that if $\log n \geq m^2(\epsilon/8, 0.99)$ then there are two distinct distributions q_1, q_2 over $[\log n]$,
 585 with disjoint support such that $\text{IPM}_{\mathcal{G}_k}(q_1, q_2) < \epsilon$. The proof is done as in claim 3.

586 Indeed, we draw IID, and uniformly, two random samples S_1 and S_2 from $\{1, \dots, \log n\}$ of size
 587 $m(\epsilon/8, 0.99)$. One can show that w.p 1/4 we have that $S_1 \cap S_2$ are distinct, also we have w.p 0.98
 588 that $\text{IPM}_{\mathcal{G}}(p_S, p_{S_1}) < \epsilon/8$ and similarly $\text{IPM}_{\mathcal{G}}(p_S, p_{S_2}) < \epsilon/8$. Taken together we obtain that with
 589 positive probability $q_1 = p_{S_1}$ and $q_2 = p_{S_2}$ have disjoint support and $\text{IPM}_{\mathcal{G}}(q_1, q_2) < \frac{\epsilon}{4}$.

590 As in theorem 5, let A be the support of q_1 and consider a distribution $p_1 = \frac{1}{k} \delta_{v_A} + (1 - \frac{1}{k}) q_1$ and
 591 similarly $p_2 = \frac{1}{k} \delta_{v_A} + (1 - \frac{1}{k}) q_2$. One can show that $\text{IPM}_{\mathcal{G}}(p_1, p_2) < \frac{\epsilon}{4}$ but the probability to draw
 592 an edge from g according to q_1 is at least 1/4, while it equals 0 if we draw edges according to p_2 .

593 To conclude, we showed that if $\log n \geq m^2(\epsilon/8, 0.99)$ then $\text{IPM}_{\mathcal{G}_k} \prec \epsilon \text{IPM}_{\mathcal{G}}$. In other words, if
 594 $\text{IPM}_{\mathcal{G}_k} \succ \epsilon \cdot \text{IPM}_{\mathcal{G}}$ then $\log n \leq m^2(\epsilon/8, 0.99)$.

$$\rho = \Omega\left(\frac{\epsilon^2}{k^2} \sqrt{\log n}\right).$$

595 C.4 Proof of theorem 7

596 **Theorem 7.** *Let $|\mathcal{V}| = n$. There exists a 2-distinguishing class \mathcal{G} , with sample complexity $m(\epsilon, \delta) =$
 597 $O(\frac{\log 1/\delta}{\epsilon^2})$ (in fact $|\mathcal{G}| = 1$) such that: For any $\epsilon > 0$ and any distinguishing class \mathcal{D} if:*

$$\text{IPM}_{\mathcal{D}} \succ \epsilon \cdot \text{IPM}_{\mathcal{G}}$$

598 *then $\text{gVC}(\mathcal{D}) = \tilde{\Omega}(\epsilon^2 \log n)$.*

599 The proof is similar to the proof of theorem 4 but we will use an improved upper bound on the size of
 600 S which we next state (see appendix C.5 for a proof):

601 **Lemma 3.** *Let \mathcal{D} be a class with $\text{gVC}(\mathcal{D}) = \rho$ over a domain S . There exists a constant $c > 0$
 602 (independent of \mathcal{D} and d) such that if $|S| > c \cdot \frac{d}{\epsilon^2} \log^2(d/\epsilon^2)$, Then there are two distributions q_1 and
 603 q_2 , supported on S such that:*

- 604 1. q_1 , and q_2 have disjoint support.
- 605 2. $\text{IPM}_{\mathcal{D}}(q_1, q_2) < \epsilon$

606 The graph g is constructed as in theorem 4. Let \mathcal{V} be a set of vertices of size $n + \log n$, let \mathcal{V}_1 be a
 607 set of size $\log n$ and we index its elements with $\{v_1, \dots, v_2, \dots, v_{\log n}\}$. We let \mathcal{V}_2 include all other
 608 elements and we index them via subsets of $[\log n]$. The graph is again constructed so that $v_A \in \mathcal{V}_2$
 609 has an edge to $v_i \in \mathcal{V}_1$ iff $i \in A$. As before, we make the graph bipartite, i.e. both \mathcal{V}_1 and \mathcal{V}_2 are
 610 independent sets.

611 Now suppose $\log n \geq c \frac{\rho}{\epsilon^2} \log^2 \frac{d}{\epsilon^2}$. By lemma 3 we have that there exists a set $A \subseteq \{1, \dots, \log n\}$, a
 612 distribution p_1 and p_2 where p_1 is supported on A and p_2 is supported on its complement so that
 613 $\text{IPM}_{\mathcal{G}}(p_1, p_2) < \epsilon$. As before we construct $q_1 = \delta v_A + (1 - \delta)p_1$ and $q_2 = \delta v_A + (1 - \delta)p_2$. One
 614 can verify that $\text{IPM}_{\mathcal{G}}(q_1, q_2) < \epsilon$ but $\text{IPM}_{\mathcal{G}_{k+1}}(q_1, q_2) > \frac{1}{2}$. Thus, if $\text{IPM}_{\mathcal{G}_k} \succ \epsilon \cdot \text{IPM}_{\mathcal{G}_{k+1}}$ then
 615 $\log n \leq c \frac{\rho}{\epsilon^2} \log^2 \frac{d}{\epsilon^2}$. In turn $d = \tilde{\Omega}(\epsilon^2 \log n)$.

616 C.5 Proof of lemma 3

617 First w.l.o.g we assume that the constant functions are in \mathcal{D} (i.e. 0 and 1).

618 We want to choose a constant c so that if $|S| \geq c \frac{2d}{\epsilon^2} \log^2 \frac{2d}{\epsilon^2}$, then we have $\frac{|S|}{\ln^2 |S|} > \frac{2\rho \log e}{\epsilon^2}$. Fix such
 619 $c > 0$, and let $\mathcal{H}_m = \{\text{sign}(\sum_{i=1}^m (2d_i(v) - 1)) : d_i \in \mathcal{D}\}$ and denote $\mathcal{H} = \mathcal{H}_{\frac{2}{\epsilon^2} \ln |S|}$. Note that

$$\begin{aligned} |\mathcal{H}| &\leq |\mathcal{D}|^{\frac{2}{\epsilon^2} \ln |S|} \\ &\leq |S|^{\frac{2\rho}{\epsilon^2} \ln |S|} && \text{Sauer's Lemma} \\ &= 2^{\frac{2\rho \log e}{\epsilon^2} \ln^2 |S|} \\ &< 2^{|S|} \end{aligned}$$

620 It thus follows that there exists $f \notin \mathcal{M}$. Let f be such and define a matrix $M = \{0, 1\}^{|S| \times |\mathcal{D}|}$ so that

$$M_{v,d} = \begin{cases} 1 & d(v) \neq f(v) \\ 0 & \text{else} \end{cases}$$

621 Now suppose that for some distribution q over S , for every d we have that $\mathbb{E}_{v \sim q}[d(v) = f(v)] < \frac{1}{2} + \frac{1}{\epsilon}$.
 622 Then, defining $q_1 = q(\cdot | f(v) = 0)$ and $q_2 = q(\cdot | f(v) = 1)$ yields the desired result. Indeed,

$$\begin{aligned}
 \sup_{d \in \mathcal{D}} |\mathbb{E}_{q_1}[d] - \mathbb{E}_{q_2}[d]| &= \sup_{d \in \mathcal{D}} 2 \left| \frac{1}{2} \mathbb{E}_{q_1}[d] - \frac{1}{2} \mathbb{E}_{q_2}[d] \right| \\
 &\geq \sup_{d \in \mathcal{D}} 2 |q(f(v) = 1) \mathbb{E}_{q_1}[d] - q(f(v) = -1) \mathbb{E}_{q_2}[d]| - 4 \max_{y \in \{1, -1\}} \left\{ \left| \frac{1}{2} - q(f(v) = y) \right| \right\} \\
 &\sup_{d \in \mathcal{D}} 2 \left| \mathbb{E}_{(v, y) \sim q} y d(v) \right| - 4\epsilon \\
 &= \sup_{d \in \mathcal{D}} 2 |1 - 2q(d(v) \neq f(v))| - 4\epsilon \\
 &\geq 8\epsilon.
 \end{aligned}$$

623 We now wish to prove that indeed, such a q exists. Suppose, otherwise: That for any distribution q
 624 over S we can find d such that $\mathbb{E}_{v \sim q}[d(v) = f(v)] > \frac{1}{2} + \frac{1}{\epsilon}$. This can be rephrased in terms of a
 625 value of a minimax game as follows:

$$\max_{q \in \Delta(S)} \min_{d \in \mathcal{D}} q^\top M_d < \frac{1}{2} - \epsilon,$$

626 Where $\Delta(S)$ denotes the set of distributions over S . It is well known ([16], thm 2), that for any game
 627 defined by any matrix M with c columns, there exists a strategy for the row player that chooses
 628 uniformly from a multiset of $\frac{\ln c}{2\epsilon^2}$ and achieves ϵ -optimality.

629 In our setting, this translate to a uniform distribution p , supported on $\frac{\ln |S|}{2\epsilon^2}$ distinguishers
 630 $\{d_1, \dots, d_{\frac{\ln |S|}{2\epsilon^2}}\}$ such that

$$\frac{2\epsilon^2}{\ln |S|} \sum_{d_i} [d_i(v) \neq f(v)] < \frac{1}{2}$$

631 , this contradicts the fact that $f \notin \mathcal{M}$.

632 We thus obtain that there exists a distribution q over S so that for every $d \in \mathcal{D}$ $\sum q(v)[d(v) \neq$
 633 $f(v)] > \frac{1}{2} - \epsilon$.

634 D Additional Proofs

635 D.1 Proof of claim 1

636 Consider the Vandermonde Matrix $V \in M_{k+1, k+1}$ given by $V_{i,j} = \left(\frac{i-1}{k}\right)^{j-1}$. Our first step will
 637 be to lower bound the smallest singular value of V . In turn, we will obtain a lower bound on the
 638 maximum value over the coordinates of the vector $V\mathbf{a}$. The proof can then be derived from the
 639 identity: $(V\mathbf{a})_i = \sum_{j=1}^{k+1} a_j \left(\frac{i-1}{k}\right)^j$.

640 Let $\lambda_1 \leq \lambda_2 \leq \dots \leq \lambda_{k+1}$ be the singular values of V . To bound the smallest singular value, λ_1 , we
 641 first observe that λ_{k+1} —the highest singular value—is bounded by $k+1$. To see that $\lambda_{k+1} \leq k+1$,
 642 observe that for any vector $\|\mathbf{a}\| \leq 1$ we have that

$$\|V\mathbf{a}\|_2 \leq k+1 \max |V_{i,j}| |a_i| \leq k+1.$$

643 Next, using the formula for the determinant of a Vandermonde matrix, and the relation $\det(V) =$
 644 $\prod \lambda_i$, we obtain:

$$\begin{aligned}
 \prod_{i=1}^{k+1} |\lambda_i| &= |\det(V)| \\
 &= \prod_{1 \leq i < j \leq k+1} \frac{|i-j|}{k} \\
 &\geq 2^{-\frac{k(k-1)}{2} \log k}
 \end{aligned}$$

645 Taken together we obtain

$$\begin{aligned}
\lambda_{\min} &\geq \frac{2^{-\frac{k(k-1)\log k}{2}}}{\prod_{i=2}^{k+1} \lambda_i} \\
&\geq \frac{2^{-\frac{k(k-1)\log k}{2}}}{\lambda_{k+1}^k} \\
&\geq 2^{-k(k-1)\log k - k\log(k+1)} \\
&= 2^{-k^2 + k\log k/k+1} \\
&\geq 2^{-2k^2}
\end{aligned}$$

646 Finally, for any polynomial $p = \sum a_i q^i$ with coefficient $|a_1|$ we have that $\|\mathbf{a}\|_2 \geq |a_1|$. We thus
647 obtain,

$$\begin{aligned}
\max_i p\left(\frac{i}{q}\right) &\geq \frac{1}{\sqrt{k+1}} \sqrt{\sum |p\left(\frac{i}{q}\right)|^2} \\
&= \frac{1}{\sqrt{k+1}} \|\mathbf{V}\mathbf{a}\|_2 \\
&\geq \frac{1}{\sqrt{k+1}} \lambda_1 \|\mathbf{a}\|_2 \\
&\geq \frac{2^{-2k^2 - 1/2\log(k+1)}}{k+1} |a_1| \\
&\geq 2^{-3k^2} |a_1|
\end{aligned}$$