003 004

010 011

012

013

014

015

016

017

018

019 020

022

#### EMBEDDING-BASED STATISTICAL INFERENCE 002 ON GENERATIVE MODELS

**Anonymous authors** 

Paper under double-blind review

## ABSTRACT

Generative models are capable of producing human-expert level content across a variety of topics and domains. As the impact of generative models grows, it is necessary to develop statistical methods to understand the population of available models. These methods are particularly important in settings where the user may not have access to information related to a model's pre-training data, weights, or other relevant model-level covariates. In this paper we extend recent results on representations of black-box generative models to model-level statistical inference tasks. We demonstrate – both theoretically and empirically – that the use of these representations are effective for multiple model-level inference tasks.

INTRODUCTION 1 021

Generative models have recently met or surpassed human-level standards on benchmarks across 023 a range of tasks (Nori et al., 2023; Katz et al., 2024; Dubey et al., 2024). While these claims warrant skepticism and further robustness evaluation (Ness et al., 2024), the impressive capabilities 025 have created a competitive environment for training state-of-the-art models and have inspired the 026 development of complementary methods to adapt models to particular use cases. For example, 027 quantizing the weights of a neural model enables a trade-off of model precision with vRAM and disk space (Gholami et al., 2022) while methods such as Low Rank Adaptation (LoRA) (Hu et al., 029 2021) and prompt-tuning (Lester et al., 2021) enable compute- and data-efficient model adaptation. Other methods such as retrieval-augmented generation (Lewis et al., 2020), model merging (Matena 031 & Raffel, 2022), constrained decoding (Hokamp & Liu, 2017), etc. (Dettmers et al., 2024; Edge 032 et al., 2024), have similarly contributed to the rapid development of a large population of diverse and accessible models. 033

034 Each model in the population has an accompanying set of covariates – scores on benchmarks, training mixture proportions, model safety scores, etc. - that are a function of the model, the training 036 set, the architecture, the retrieval database, or derivatives thereof. For a given model, the covariate 037 of interest may not be available to the user or it may be too expensive to calculate. For example, it 038 may not be known if a proprietary model has been trained on copyrighted data. Or, it may be too resource intensive to directly evaluate how toxic every model is. Methods for predicting model-level covariates are necessary in these settings, and others, to fully understand the behavior and properties 040 of a model. 041

042 In this paper we extend recent theoretical results for vector representations of generative models 043 (Acharyya et al., 2024) to model-level statistical inference settings. Our results show that the em-044 beddings of the responses from a collection of generative models can be used for consistent inference for a wide class of inference problems. We demonstrate the effectiveness of the representations for 045 three downstream inference tasks, including predicting the presence of sensitive information in a 046 model's training mixture and predicting a model's safety. We include empirical investigations of 047 performance sensitivity to hyperparameters required to generate the model-level vector representa-048 tions. 049

**Contribution.** Our contribution is the theoretical and empirical validation for using vector repre-051 sentations of black-box generative models for model-level inference tasks. The representations that we study herein are based on the embeddings of their responses. The theoretical results apply to any 052 generic, well-behaved embedding function. Given the representations, any standard vector-based method can be used for inference on the models.

# 054 1.1 BACKGROUND & RELATED WORK

056 Our work is an extension of recent theoretical and empirical results on embedding-based represen-057 tations of generative models (Acharyya et al., 2024), which itself is a continuation of a long-line 058 of embedding-based investigations of the inputs and outputs of generative models (Mikolov, 2013; Reimers, 2019; Neelakantan et al., 2022; Patil et al., 2023). Of particular relevance is recent empirical work defining the data kernel (Duderstadt et al., 2023) and investigations into its ability to track 060 the dynamics of interacting models (Helm et al., 2024), in which the experiments demonstrate the 061 ability to parlay the embeddings of a collection of inputs or outputs into useful vector representations 062 of the generative models themselves in both white-box and black-box settings. The results herein 063 further theoretically and empirically validate these findings in the context of statistical inference. 064

Our work is also related to using embedding-based techniques for inference on complicated objects such as entire mouse connectomes (Wang et al., 2020), physiological data (Chen et al., 2022), and classification distributions (Helm et al., 2021). In each of these settings, the authors define a pairwise distance matrix on the objects and apply multi-dimensionsal scaling to obtain vector representations of each of the objects. Once there is one vector representation per object, standard inference methods for the specific task can be used. The method we study herein follows this general formula, with the additional complication that the objects are random mappings.

Lastly, our work is a part of the relatively new literature on inference on generative models. For
example, FlashHELM (Perlitz et al., 2023) uses the score on a subset of a benchmark such as HELM
(Liang et al., 2022) to quickly identify where a model fits on a leaderboard. More recent work
proposes scaling laws to predict the performance of base models on a suite of benchmarks as a
function of training FLOPs (Ruan et al., 2024). The method proposed herein does not assume access
to a scoring function nor does it assume access to a function of the model's weights at inference time.
Perhaps most importantly, our setting and method can be applied to general model-level prediction.

079 080

081

082

098 099

## 2 THE DATA KERNEL PERSPECTIVE SPACE

Before we present our results related to statistical inference on black-box generative models, we
 first describe how to obtain finite-dimensional vector representations of the models. The particular
 method for obtaining vector representations for generative models that we study herein has previously been discussed in Helm et al. (2024) and Acharyya et al. (2024).

Let  $f : \mathcal{Q} \to \mathcal{X}$  be a black-box model from a query space  $\mathcal{Q}$  to an output space  $\mathcal{X}$ . In our setting there are *n* models  $f_1, \ldots, f_n$  and *m* queries  $q_1, \ldots, q_m$  with  $q_j \in \mathcal{Q}$ . We assume that each model responds to every query *r* independent times and let  $f_i(q_j)_k$  be the *k*-th response to  $q_j$  from  $f_i$ .

We let  $g: \mathcal{X} \to \mathbb{R}^p$  be an embedding function that maps a model response to a *p*-dimensional realvalued vector. Further, we let  $x_{ijk} = g(f_i(q_j)_k)$  denote the embedding of  $f_i(q_j)_k$  and  $F_{ij}$  denote the distribution of  $x_{ijk}$  on  $\mathbb{R}^p$ . That is,  $x_{ij1}, \ldots x_{ijr} \sim^{iid} F_{ij}$ .

We let  $\bar{x}_{ij} = \frac{1}{r} \sum_{k=1}^{r} x_{ijk}$  be the empirical average embedded response of  $f_i$  to  $q_j$  and  $\bar{X}_i \in \mathbb{R}^{m \times p}$ denote the matrix whose *j*-th row is  $\bar{x}_{ij}$ . We view  $\bar{X}_i$  as a matrix representation of model  $f_i$  with respect to  $\{q_1, \ldots, q_m\}$  and *g*. We define *D* as the  $n \times n$  pairwise distance matrix with entries

$$D_{ii'} = \frac{1}{m} \big| \big| \bar{X}_i - \bar{X}_{i'} \big| \big|_F.$$
(1)

The matrix entry  $D_{ii'}$  captures the difference in empirical average model behavior between  $f_i$  and  $f_{i'}$  with respect to  $\{q_1, \ldots, q_m\}$  and g.

With the form described by Eq. (1), D is a rescaled Euclidean distance matrix. Hence, multidimensional scaling (MDS) of D yields d-dimensional Euclidean representations of the matrices  $\bar{X}_i$  (and thereby of the models  $f_i$ ) with respect to  $\{q_1, \ldots, q_m\}$  and g (Torgerson, 1952). Letting  $\hat{\psi} := \text{MDS}(D) \in \mathbb{R}^{n \times d}$ , we refer to  $\hat{\psi}$  as the *data kernel perspective space* (DKPS). We emphasize that the *i*-th row of the DKPS –  $\hat{\psi}_i$  – is a d-dimensional vector representation of the generative model  $f_i$ .

## 2.1 ANALYTICAL PROPERTIES OF THE DATA KERNEL PERSPECTIVE SPACE

110 While the DKPS representations of the models are Euclidean objects, it is not possible to comment 111 on their properties without imposing constraints on the  $F_{ij}$ . We let  $\mu_i \in \mathbb{R}^{m \times p}$  be the population 112 counterpart of  $\bar{X}_i$  whose *j*-th row is  $\mathbb{E}_{x \sim F_{ij}}(x)$ . Similarly, we let  $\Delta$  be the population counterpart 113 of *D* whose (i, i')-th element is  $\Delta_{ii'} = \frac{1}{m} ||\mu_i - \mu_{i'}||_F$ . We assume  $\Delta_{ii'}^* = \lim \frac{1}{m} ||\mu_i - \mu_{i'}||_F$  for 114 every pair (i, i') as  $m, r \to \infty$  and that *g* is bounded.<sup>1</sup>

In our setting, under appropriate assumptions described in (Acharyya et al., 2024),

117

127

128

136 137 138

144 145

156 157 158

161

 $D_{ii'} = \frac{1}{m} \|\bar{X}_i - \bar{X}_{i'}\| \to \Delta^*_{ii'}$ (2)

with high probability as  $m, r \to \infty$  for all  $(i, i') \in \{1, ..., n\} \times \{1, ..., n\}$ . Further, there exists  $\psi := \text{MDS}(\Delta^*) \in \mathbb{R}^{m \times d}$  such that  $\widehat{\psi} \to \psi$ . That is, the configuration  $\psi$  is the population counterpart of  $\widehat{\psi}$ .  $\psi$  captures the true geometry of the mean discrepancies of the model responses with respect to the queries  $\{q_1, ..., q_m\}$ . Importantly, in settings where the models are not fixed and are instead assumed to be *i.i.d.* realizations from a model distribution  $P_{model}$ , the distance matrix Dconverges to  $\Delta^*$  and, under technical assumptions, there exists  $\psi$  such that  $\widehat{\psi} \to \psi$  as  $m, r \to \infty$ for all n (Acharyya et al., 2024).

#### 2.2 STATISTICAL INFERENCE IN THE DATA KERNEL PERSPECTIVE SPACE

We now extend the consistency of  $\hat{\psi}$  to  $\psi$  to model-level inference. Consider the classical statistical learning problem (Hastie et al., 2009, Chapter 2) in the context of a collection of generative models: given training data  $\mathcal{T}_n = \{(f_1, y_1), \dots, (f_n, y_n)\}$  assumed to be *i.i.d.* realizations from the joint distribution  $P_{fY}$ , choose the decision function  $h : \mathcal{F} \to \mathcal{Y}$  that minimizes the expected value of a loss function  $\ell$  with respect to  $P_{fY}$  within a class of decision functions  $\mathcal{H}$  for a test observation fassumed to be an independent realization from the marginal distribution  $P_f$ . Or, with  $\mathcal{R}_{\ell}(P_{fY}, h) :=$  $\mathbb{E}_{P_{fY}} [\ell(h(f), y)]$ , select  $h^*$  such that

$$h^* \in \operatorname*{argmin}_{h \in \mathcal{H}} \mathcal{R}_{\ell}(P_{fY}, h)$$

139 We let  $\mathcal{R}^*_{\ell}(P_{fY}, \mathcal{H}) = \mathcal{R}_{\ell}(P_{fY}, h^*)$  denote the expected loss (or "risk") of  $h^*$ .

The joint distribution  $P_{fY}$  is often unavailable and the decision function is selected based on the training data. We let  $h(\cdot; \mathcal{T}_n)$  denote such a decision function and say the sequence of decision functions  $(h(\cdot; \mathcal{T}_1), \dots, h(\cdot; \mathcal{T}_n))$  is consistent for  $P_{fY}$  with respect to  $\mathcal{H}$  if

 $Pr(\left|\mathcal{R}_{\ell}(P_{fY}, h(\cdot; \mathcal{T}_n)) - \mathcal{R}_{\ell}^*(P_{fY}, \mathcal{H})\right| > \epsilon) \to 0$ 

146 as  $n \to \infty$ .

147 In the black-box setting we do not have direct access to useful representations of the models. Instead, 148 we can use the representations  $\{\psi_1, \ldots, \psi_n\}$  as proxies for the models  $\{f_1, \ldots, f_n\}$ . Hence,  $\mathcal{T}_n =$ 149  $\{(\psi_1, y_1), \ldots, (\psi_n, y_n)\}$  where  $(\psi_i, y_i) \sim^{iid} P_{\psi Y}$  and our goal is to choose a decision function  $h^*$ 150 which minimizes the average loss  $\mathcal{R}_{\ell}(P_{\psi Y}, h) := \mathbb{E}_{P_{\psi Y}} [\ell(h(\psi), y)]$  within an appropriately defined 151  $\mathcal{H}$ . Note that the true  $\{\psi_1, \ldots, \psi_n\}$  are not known *a priori* and must be estimated via  $\{\widehat{\psi}_1, \ldots, \widehat{\psi}_n\}$ . 152 We let  $\widehat{\mathcal{T}}_n$  be the training data where  $\psi_i$  is replaced with  $\widehat{\psi}_i$ .

Two important extensions of the DKPS consistency results in the context of the classical statistical learning problem are the following theorems:

**Theorem 1.** Under technical assumptions described in Appendix A,

$$\mathcal{R}_{\ell}(P_{\psi Y}, h(\cdot; \widehat{\mathcal{T}}_n)) \to \mathcal{R}_{\ell}(P_{\psi Y}, h(\cdot; \mathcal{T}_n))$$

159  $as m, r \to \infty$ , for every n.

<sup>&</sup>lt;sup>1</sup>Boundedness is typically satisfied in practice for well-defined g. For language models, an element of  $\mathcal{X}$  is a finite sequence from a finite vocabulary; for text-to-image models  $\mathcal{X} = \{0, \dots, 255\}^3$  is finite.

That is, the risk of a decision function based on  $\hat{\mathcal{T}}_n = \{(\hat{\psi}_1, y_1), \dots, (\hat{\psi}_n, y_n)\}$  converges to the risk of the decision function based on the true-but-unknown  $\mathcal{T}_n$ . For situations where  $\{\psi_1, \dots, \psi_n\}$  are good proxies, Theorem 1 states that increasing the number of queries and/or number of replicates per query will improve the performance of decision functions trained on  $\hat{\mathcal{T}}_n$ .

**Theorem 2.** Under technical assumptions described in Appendix A, if  $(h(\cdot; \mathcal{T}_1), \ldots, h(\cdot; \mathcal{T}_n))$ is consistent for  $P_{\psi Y}$  with respect to  $\mathcal{H}$ , then  $(h(\cdot; \hat{\mathcal{T}}_1), \ldots, h(\cdot; \hat{\mathcal{T}}_n))$  is consistent for  $P_{\psi Y}$  with respect to  $\mathcal{H}$  as  $n, m, r \to \infty$ .

Theorem 2 states that if the sequence of decision functions learned from  $\mathcal{T}_n$  is consistent then the sequence of decision functions learned from the analogous  $\hat{\mathcal{T}}_n$  is also consistent. This result suggests that inference performance should improve with more training data. While a logical follow-on to Theorem 1 and a well-understood principle in classical statistical learning settings, the consistency of the DKPS in the setting where the number of models grows is a non-trivial extension of the fixed n case.

- The proofs of Theorems 1 and 2 are provided in Appendix A.
- 177 178

2.2.1 NON-STANDARD CONSIDERATIONS

Given that our analysis is based on estimates  $\{\hat{\psi}_1, \dots, \hat{\psi}_n\}$  of proxies  $\{\psi_1, \dots, \psi_n\}$  of the objects of interest  $\{f_1, \dots, f_n\}$ , it is important to note some non-standard theoretical considerations. To facilitate our discussion we let  $\mathcal{R}^*_{\ell}(P_{fY}) := \inf_{\{h: \mathcal{F} \to \mathcal{Y}\}} \mathcal{R}_{\ell}(P_{fY}, h)$  be the Bayes optimal risk of  $P_{fY}$ . Further, we let  $P_{query}$  be a distribution on queries.

For example, while our theoretical results provide insights as to how performance will be affected by increasing n, m, and r, our results do not comment on the magnitude of  $|\mathcal{R}_{\ell}^{*}(P_{\psi Y}) - \mathcal{R}_{\ell}^{*}(P_{fY})|$ . This quantity is a measure of how good of a proxy  $\{\psi_{1}, \ldots, \psi_{n}\}$  is for  $\{f_{1}, \ldots, f_{n}\}$ . In particular, if  $|\mathcal{R}_{\ell}^{*}(P_{\psi Y}) - \mathcal{R}_{\ell}^{*}(P_{fY})| = 0$  then the  $\psi_{i}$  are perfect proxies of the  $f_{i}$  with respect to y (Devroye et al., 2013, Chapter 32). We expect it to be challenging to theoretically understand  $|\mathcal{R}_{\ell}^{*}(P_{\psi Y}) - \mathcal{R}_{\ell}^{*}(P_{fY})|$ in general.

Similarly, let  $q_1, \ldots, q_m \sim^{iid} P_{query}$  and  $q'_1, \ldots, q'_m \sim^{iid} P'_{query}$  with  $P_{query} \neq P'_{query}$ . The proxies of the models induced by these two query sets  $-\{\psi_1, \ldots, \psi_n\}$  and  $\{\psi'_1, \ldots, \psi'_n\}$ , respectively – will be of different quality. Without an initial understanding of the models and their relationship with the covariates, we expect that it will be impossible to know which query distribution is preferred *a priori*. We highlight the importance of  $P_{query}$  in an experimental setting below.

Lastly, we note that representations of the models that we analyze herein capture the intrinsic geometry of the mean discrepancies of the models for the queries  $\{q_1, \ldots, q_n\}$ . For covariates that cannot be described as a function of the mean discrepancies of the models, it is unclear how to interpret Theorems 1 and 2. We discuss potential alternatives to the representations that we study in the discussion.

201 202 2.3 AN ILLUSTRATIVE EXAMPLE – "WAS RA FISHER GREAT?"

The first model-level inference task we study is a toy example where the task is to predict the probability that a model will respond "yes" to the question "Was RA Fisher great?". The question is chosen due to its subjectiveness – there is no correct answer to the question of a person's greatness – and its duality – Ronald A. Fisher is considered one of the most influential statisticians in history and is considered an advocate of the 20th century Eugenics movement (Bodmer et al., 2021).

208 We use a 4-bit version of Meta's LLaMA-2-7B-Chat (Touvron et al., 2023) as a base model and 209 consider a collection of models parameterized by fixed context augmentations. Each augmentation 210  $a_i$  contains information related to RA Fisher's statistical achievements (i.e.,  $a_i =$ "RA Fisher pio-211 neered the principles of the design of experiments") or to his involvement in the eugenics movement 212 (i.e.,  $a_i =$  "RA Fisher's view on eugenics were primarily based on anecdotes and prejudice.") and is 213 prepended to every query. The covariate corresponding to a given model is calculated by prompting the base model with the appropriately formatted prompt "Give a precise answer to the question 214 based on the context. Don't be verbose. The answer should be either a yes or a no.  $a_i$ . Was RA 215 Fisher great?" until there are 100 valid responses. We let  $y_i$  be the average number of "yes" es.



Figure 1: Left. The 2-d Data Kernel Perspective Space (DKPS) and covariate surface for a collection of 550 models parameterized by fixed augmentations. **Right.** The performance of the 1-nearest neighbor regressor in DKPS for predicting the probability that an unlabeled model responds "yes" to "Was RA Fisher great?".

238 To induce a DKPS for this task, we consider queries sampled from OpenAI's ChatGPT with the 239 prompt "Provide 100 questions related to RA Fisher.". For a given query  $q_j$  we prompt the base 240 model with the appropriately formatted prompt " $a_i q_j$ " and fix r = 1. The left figure of Figure 1 is a 3-d figure where the first two dimensions are the DKPS of the n = 550 (275 statistics augmen-241 242 tations, 275 eugenics augmentations) models induced with m = 100 queries. Model responses are embedded by averaging the per-token last layer activation of the base model. The third dimension 243 is an interpolated  $y_i$  surface with a linear kernel (Du Toit, 2008). The first dimension of the DKPS 244 is clearly capable of distinguishing between models adorned with "statistics" augmentations and 245 models adorned with "eugenics" augmentations. The shape of the interpolated covariate surface is 246 highly correlated with this feature. A description of the augmentations that parameterize the models 247 and the queries used to induce the DKPS is provided in Appendix B.1. 248

The right figure of Figure 1 shows the performance of a 1-nearest neighbor (1-NN) regressor in 249 DKPS for a varying number of labeled models and a varying number of queries. The DKPS is 250 induced with n models and the regressor is trained with n-1 of the model-level covariates. The 251 mean squared error reported is the error of the 1-NN regressor for predicting the "left out" model's 252 covariate ( $\pm 1$  S.E.). As Theorems 1 and 2 suggest, the performance of the regressor is dependent 253 on both the number of models and the number of queries: the more models and the more queries 254 the better. The scale of the impact of more models and more queries, however, depends on its 255 counterpart. The amount of models does not have a large impact on predictive performance if the 256 number of queries is small. The amount of queries has a large impact on performance regardless of 257 the number of models.

258

232

233

234

235 236 237

259 260

## 3 EXPERIMENTS

261 262

263

We next consider two experiments with more realistic model-level covariates: predicting whether or not a model has had access to sensitive information and predicting model safety. For all experiments we fix r = 1 as suggested by the empirical rates of convergence of the perspectives described in Acharyya et al. (2024). We use the MDS implementation from Graspologic (Chung et al., 2019) throughout. We use the profile likelihood of the singular values of D to determine the dimensionality of the DKPS (Zhu & Ghodsi, 2006) and note that this may be larger than two. We show only the first two dimensions for visualization purposes.



Figure 2: Left. The 2-d data kernel perspective space (DKPS) of 50 fine-tuned models – 25 with "sensitive" data in the fine-tuning data mixture (red), 25 with none (black) – induced by an evaluation set containing 10 prompts relevant to the sensitive data. For models trained on sensitive data, color intensity correlates with amount of sensitive data in the training mixture. Center. The 2-d DKPS of the models induced by a set of 10 prompts "orthogonal" to the difference between models with sensitive data in their fine-tuning data mixture and models with no sensitive data in their fine-tuning data mixture. Right. Classification performance as a function of number of labeled models and size of evaluation set for both sensitive and orthogonal evaluation sets.

292

282

283

284

285

286

287

288

#### 3.1 HAS A MODEL SEEN SENSITIVE INFORMATION?

293 Modern language models are trained with trillions of tokens of text (Touvron et al., 2023). For proprietary models such as OpenAI's GPT series or Anthropic's Claude series, the exact sources 295 of the training mixtures are unknown and – given the models' propensities to produce content that 296 is strikingly similar to copyrighted content (Henderson et al., 2023) - its curation and use is eth-297 ically questionable (Lemley & Casey, 2020). Further, the training mixture of some models may 298 include sensitive informative such as personal information, trade secrets, or government-classified 299 information that should never be presented to the end user. Developing classifiers to identify models that are either trained on sensitive or copyrighted information or are likely to produce sensitive or 300 copyrighted information is thus paramount to uphold the rights of the stakeholders of the original 301 content. There has been recent work on this topic in settings with access to model weights or token 302 likelihoods (see, e.g., (Duderstadt et al., 2023; Shi et al., 2023)). 303

304 To investigate the utility of DKPS for this purpose, we again use a 4-bit version of LLaMA-2-7B-Chat as a base model and train 50 different LoRA adapters with different sub-305 sets of the Yahoo! Answers (YA) dataset (Zhang et al., 2015). The YA dataset consists of data from 306 10 topics. We consider data from the topic "Politics & Government" to be "sensitive" information, 307 data from the topics "Society & Culture", "Science & Mathematics", "Health", "Education & Ref-308 erence", "Computers & Internet", and "Sports" to be "not-sensitive", and data from the remaining 309 topics to be "orthogonal". We trained each of the 50 adapters with 500 question-answer pairs for 310 3 epochs with a learning rate of  $5 \times 10^{-5}$  and a batch size of 8. Each adapter is rank 8, has a 311 scaling factor of 32, targets all attention layers, does not have bias terms, and has a dropout proba-312 bility of 0.05 when training. For 25 of the adapters, the adapter training mixture consisted wholly 313 of randomly selected not-sensitive data. For the remaining 25 adapters, the adapter training mixture 314 consisted of  $p_i$  randomly selected sensitive data and  $500 - p_i$  randomly selected not-sensitive data. 315 We let  $y_i$  be the indicator of whether or not the adapter's training mixture contained any sensitive data. 316

We study classification of the models in two different DKPS: one induced by a set of randomly selected sensitive queries, one induced by a set of randomly selected orthogonal queries. Both DKPS use the open source embedding model nomic-embed-v1.5 (Nussbaum et al., 2024). A 2-d DKPS induced by m = 10 randomly selected sensitive queries is shown on the left of Figure 2. In this space the models are separated by their label and the models that have seen more sensitive information are generally farther from the class-boundary. A 2-d DKPS induced by m = 10 orthogonal queries is shown in the center of Figure 2. Here, by contrast, the models are not easily separable by their label.

324 The right figure of Figure 2 shows the classification performance of Fisher's Linear Discriminant 325 trained on varying amounts of DKPS representations of models for the two query distributions. For 326 a given n and m, we induce the DKPS with all models and a randomly selected query set. We report 327 the expected risk of the classifier trained on a random subset of the models for the remaining models. 328 We observe similar phenomena to the "Was RA Fisher great?" experiment in that both the amount of models and the amount of queries impact performance. For a fixed m, the expected risk curves also 329 highlight the observed difference in separability of the models in the two DKPS, with the expected 330 risk with sensitive queries being significantly lower than the expected risk with orthogonal queries. 331 Indeed, the expected risk with m = 10 sensitive queries is similar to the expected risk with m = 50332 orthogonal queries. 333

Lastly, we highlight the 1-d projection learned by Fisher's linear discriminant for m = 10for both DKPS in Figure 3. As can be seen in the top figure of Figure 3, the projection of the models is correlated with the proportion of sensitive data that the model has had

- access to when the DKPS is induced with sen-337 sitive queries. While the linear goodness-of-fit 338 is not large ( $R^2 = 0.37$ ), the correlation is sta-339 tistically significant per the hypothesis test us-340 ing Kendall's rank correlation coefficient ( $\tau =$ 341 0.42, p < 0.01). The projection of the mod-342 els when the DKPS is induced with orthogonal 343 queries has a line of best fit with a negligible 344 slope and a much smaller Kendall's rank cor-345 relation coefficient ( $\tau = 0.08$ ). The second of which results in a p value of 0.57 – meaning we 346 fail to reject the null that the amount of sensi-347 tive information the model has had access to is 348 correlated with the learned 1-d projection. 349
- 350 Throughout this experiment we use the term 351 "orthogonal" for queries from topics that are not used when fine-tuning the models under 352 study. The term is chosen because, naïvely, 353 queries from these topics should not elicit dif-354 ferent responses from models trained on sen-355 sitive data and models trained on not-sensitive 356 data. In reality, the sensitive data and the not-357 sensitive data have different underlying token 358 distributions and this difference will cause sys-359 tematic differences in model responses after 360 fine-tuning even for topics that are irrelevant 361 a priori. Further, it is likely that the docu-362 ments in the "orthogonal" topics share some content commonalities with documents in the 363 sensitive and not-sensitive topics. We see this 364



Figure 3: **Top.** The 1-d FLD projection of the models from a DKPS induced by queries from the sensitive topic versus the amount of sensitive data the adapter had access to during training. **Bottom.** The same but for a DKPS induced by queries from orthogonal topics.

phenomenon in the classification results in the right figure of Figure 2 where the linear classifier is
 able to perform better than chance with enough "orthogonal" queries.

367 368

369

#### 3.2 HOW SAFE IS A MODEL?

Model safety is one of the biggest concerns when deploying a language model in production. An unsafe model is prone to propagating harmful stereotypes (Ferrara, 2024), using toxic language (Wen et al., 2023), and misunderstanding the user's intent (Ji et al., 2023) – all of which can adversely affect the user and their experience. Hence, developing techniques to understand how unsafe a model is an important aspect of the model-production pipeline. As with predicting if a model has had access to sensitive information above, we investigate using DKPS to predict model safety through the lens of model toxicity and model bias.

We consider a collection of 58 models. Each model in the collection is a base model, a fine-tuned version of a base model, or the result of weight-merging various other models in the collection. We

383

387

391 392

394

396

397



Figure 4: Left. A graph where each node is a model and an edge between two models exists if model i is fine-tuned from model i' or if model i's weights were used in a model-merge that 393 resulted in model *i*, etc. **Right, top.** The two-dimensional data kernel perspective spaces (DKPS) corresponding to the toxicity and bias prediction tasks. Dot size is proportional to model toxicity or 395 bias. **Right, bottom.** Average relative performance of three regression techniques across all models in the model graph. Local predictions in DKPS are more effective than both global predictions and local predictions in model relationship space.

398 399 400

401 view this collection of models as a graph where each model is a node and an undirected edge exists between nodes if one of the models is a fine-tuned version of the other or if one of the models is 402 the result of a model merge including the other. The graph representing the collection of models is 403 shown on the left of Figure 4. The list of models under study is provided in Appendix B.2. 404

405 For each model in the collection we consider two covariates: model tox-406 icity and model bias. To determine а model's toxicity, we prompt each 407 model with а collection of queries from the Real Toxicity Prompts (RTP)

(Gehman et al., 2020) dataset and subsequently evalu-408 ate the toxicity of each response with the neural model 409 roberta-hate-speech-dynabench-r4 (Vidgen 410 et al., 2021). The model-level toxicity is simply the av-411 erage response toxicity. An analagous process is used to 412 define a model's bias with the dataset Bias in Open-ended 413 Language Generation Dataset (BOLD) (Dhamala et al., 414 2021) and the regard model (Sheng et al., 2019). 415

- We induce the perspective spaces for the toxicity and 416 bias tasks with randomly sampled prompts from RTP 417 and BOLD, respectively, and the embedding model 418 nomic-embed-v1.5. The 2-d DKPS - induced with 419 m = 2000 queries – for both tasks is shown in the top 420 right of Figure 4. The size of the dot is correlated with 421 the model's covariate. We highlight an example unla-422 beled model (green) and its corresponding neighbor in DKPS (red) and neighbors in graph-space (blue) in Fig-423 ure 4. Importantly, the relative position of a model in the 424 respective DKPS is predictive of the model's toxicity and 425 the model's bias. 426
- 427 We quantify this observation by evaluating regressors for 428 predicting the model-level toxicity and bias of an unla-429 beled model. We consider three different regressors for these tasks. The first is a constant equal to the average 430 covariate of the labeled models (i.e., the "global mean"). 431



Figure 5: The relative time improvement (larger is better) when using local predictions in DKPS instead of calculating the ground-truth model-level covariate using HuggingFace's API.

The second uses the average covariate of models who share an edge with the unlabaled model

432 (i.e., "1-NN (graph)"). The third uses the covariate of the nearest neighbor in DKPS (i.e., "1-NN 433 (DKPS)"). For 1-NN (DKPS) we consider varying amounts of randomly sampled queries from RTP 434 and BOLD to induce the DKPS. We use the global mean as a standard and report the relative ab-435 solute error of the three methods in the bottom right of Figure 4. The reported performance of the 436 1-NN (DKPS) regressor is the average of the smaller of 200 and 2000/m random samples of m queries. The reported relative absolute error is the average across all models in the model graph. 437 Notably, given enough queries local predictions in DKPS outperform local predictions in the graph 438 space and predictions using the global mean. 439

440 In addition to reporting the relative performance, we report the time it takes to use the DKPS machin-441 ery to predict Mistralv1.0's (Jiang et al., 2023) toxicity and bias relative to using the evaluation 442 model through HuggingFace's API (Wolf et al., 2019). The time we report for the HuggingFace API is the time it takes to calculate the "ground truth" used to calculate the performance of the regressors 443 above and is the total time required for Mistralv1.0 to produce responses and for the evaluation 444 model to produce scores for 2000 queries. The time we report for 1-NN (DKPS) includes the time 445 it takes for Mistralv1.0 to produce m responses, the time it takes nomic-embed-v1.5 to 446 embed the responses, the time it takes to induce the DKPS, and the time it takes to train and use 447 the nearest neighbor regressor. It does not include the time it takes to produce and evaluate the 448 responses of the other models in the collection. The relative efficiency of DKPS, as seen in Figure 449 5, is approximately 1/m. This relationship will hold for all model-level inference tasks where the 450 covariate is proportional to the sum of a function of individual responses. 451

452 453

454

### 4 DISCUSSION

We have demonstrated – both theoretically and empirically – that embedding-based representations of generative models can be used for various model-level inference tasks. While the results we presented show the potential of our approach, there are choices throughout the collection-of-modelsto-covariate-prediction pipeline that can affect performance and implementation practicality.

As mentioned in Section 2.2, one major decision is the query distribution  $P_{query}$  (or, more practi-459 cally, the set  $\{q_1, \ldots, q_m\}$ ). In particular, the representations of the models that the query set induces 460 may or may not be relevant to a particular inference task. We demonstrate this phenomenon in Sec-461 tion 3.1 where the queries from the "sensitive-only" query distribution can induce representations of 462 similar discriminative ability as queries from the "orthogonal-only" query distribution with 1/5 of 463 the queries. We expect a similar but less dramatic effect within the distributions of "sensitive-only" 464 queries and anticipate that curating an "optimal" set of queries for a given g and for a fixed m will 465 soon be a highly active research area. 466

Another decision is the distance function used to define D. The MDS of the distance matrix studied 467 herein (Eq. (1)) produces representations of the models that are consistent for objects that cap-468 ture the true mean discrepancy geometry of the model responses. For model-level covariates that 469 cannot naturally be described as a function of mean discrepancy geometry, we do not expect that 470 information-theoretic optimal performance when using  $\psi$  is possible for any n, m, and r. Instead, 471 for proxies of the models to minimize information loss, it is necessary to replace the Frobenius 472 norm of the differences of the average embedded response with a more expressive distance such as 473 an extension of a distance defined directly on the cumulative distributions of responses or with a 474 task-specific distance (Helm et al., 2020). Related theoretical work (Tang et al., 2013) suggests that 475 our results can be extended to more expressive distances. We do not expect the naïve replacement of 476 the Frobenius norm with a more expressive distance function to be universally better for model-level 477 inference tasks in practice as there are likely computational and query set quality trade-offs to consider. As with the active curation of an optimal query set, we expect these trade-offs to be important 478 future research topics. 479

In the experiments above we have access to n models and n' < n corresponding model-level covariates. We induce a DKPS with the entire collection of models and use the n' labeled models to predict a label for the remaining n - n' models. In practice an unlabeled model may not be available when inducing the DKPS or, if n is large, it may be too expensive to induce a DKPS whenever a prediction for a new unlabeled model is required. Out-of-sample techniques (Bengio et al., 2003; Trosset & Priebe, 2008) can be used to meet these imposed constraints at the cost of a slight degradation in representation quality and, hence, inference performance.

- Implementation of inference in DKPS in practice will require an upfront, one-time cost of generating
   responses from a subset of models in the collection and scoring their outputs. Since this is required
   to compare the models with respect to the covariates anyway we followed the timing paradigm
   presented in (Perlitz et al., 2023) and did not include this cost when comparing the methods in
   Figure 5. Once the models in the initial subset are scored we expect the relative time efficiency (and
   implied relative computational efficiency) of inference in DKPS to be worthwhile to practitioners.
- We fixed r = 1 and let m grow throughout our experiments. For Theorems 1 and 2 to hold, both m and r must grow. In practice, the trade-off between getting responses for more queries or getting more responses for the same queries depends on the distributions  $F_{ij}$ . For example, if  $F_{ij}$  is a point mass then r > 1 for  $q_j$  is unnecessary. Conversely, if  $F_{ij}$  is a complicated distribution on  $\mathbb{R}^p$  then more responses are necessary to properly estimate it and, hence, to properly capture the difference between average model responses.
- 498 499

517

522

#### References

- Aranyak Acharyya, Michael W. Trosset, Carey E. Priebe, and Hayden S. Helm. Consistent estimation of generative model representations in the data kernel perspective space, 2024. URL https://arxiv.org/abs/2409.17308.
- Yoshua Bengio, Jean-françcois Paiement, Pascal Vincent, Olivier Delalleau, Nicolas Roux, and Marie Ouimet. Out-of-sample extensions for lle, isomap, mds, eigenmaps, and spectral clustering. *Advances in neural information processing systems*, 16, 2003.
- Walter Bodmer, RA Bailey, Brian Charlesworth, Adam Eyre-Walker, Vernon Farewell, Andrew
   Mead, and Stephen Senn. The outstanding scientist, ra fisher: his views on eugenics and race.
   *Heredity*, 126(4):565–576, 2021.
- Guodong Chen, Hayden S Helm, Kate Lytvynets, Weiwei Yang, and Carey E Priebe. Mental state classification using multi-graph features. *Frontiers in Human Neuroscience*, 16:930291, 2022.
- Jaewon Chung, Benjamin D. Pedigo, Eric W. Bridgeford, Bijan K. Varjavand, Hayden S. Helm, and
   Joshua T. Vogelstein. Graspy: Graph statistics in python. *Journal of Machine Learning Research*,
   20(158):1–7, 2019. URL http://jmlr.org/papers/v20/19-490.html.
- Tim Dettmers, Artidoro Pagnoni, Ari Holtzman, and Luke Zettlemoyer. Qlora: Efficient finetuning of quantized llms. *Advances in Neural Information Processing Systems*, 36, 2024.
- Luc Devroye, László Györfi, and Gábor Lugosi. A probabilistic theory of pattern recognition, volume 31. Springer Science & Business Media, 2013.
- Jwala Dhamala, Tony Sun, Varun Kumar, Satyapriya Krishna, Yada Pruksachatkun, Kai-Wei Chang, and Rahul Gupta. Bold: Dataset and metrics for measuring biases in open-ended language generation. In *Proceedings of the 2021 ACM Conference on Fairness, Accountability, and Transparency*, FAccT '21, pp. 862–872, New York, NY, USA, 2021. Association for Computing Machinery. ISBN 9781450383097. doi: 10.1145/3442188.3445924. URL https://doi.org/10.1145/3442188.3445924.
- Wilna Du Toit. *Radial basis function interpolation*. PhD thesis, Stellenbosch: Stellenbosch University, 2008.
- Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha
  Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. The Ilama 3 herd of models. *arXiv preprint arXiv:2407.21783*, 2024.
- Brandon Duderstadt, Hayden S Helm, and Carey E Priebe. Comparing foundation models using data kernels. *arXiv preprint arXiv:2305.05126*, 2023.
- 538 Darren Edge, Ha Trinh, Newman Cheng, Joshua Bradley, Alex Chao, Apurva Mody, Steven Truitt,
  539 and Jonathan Larson. From local to global: A graph rag approach to query-focused summarization. *arXiv preprint arXiv:2404.16130*, 2024.

540 541 542	Emilio Ferrara. Fairness and bias in artificial intelligence: A brief survey of sources, impacts, and mitigation strategies. <i>Sci</i> , 6(1), 2024. ISSN 2413-4155. doi: 10.3390/sci6010003. URL https://www.mdpi.com/2413-4155/6/1/3.			
543 544 545 546	Samuel Gehman, Suchin Gururangan, Maarten Sap, Yejin Choi, and Noah A Smith. Real- toxicityprompts: Evaluating neural toxic degeneration in language models. <i>arXiv preprint</i> <i>arXiv:2009.11462</i> , 2020.			
547 548 549	Amir Gholami, Sehoon Kim, Zhen Dong, Zhewei Yao, Michael W Mahoney, and Kurt Keutzer. A survey of quantization methods for efficient neural network inference. In <i>Low-Power Computer Vision</i> , pp. 291–326. Chapman and Hall/CRC, 2022.			
550 551 552	Trevor Hastie, Robert Tibshirani, Jerome H Friedman, and Jerome H Friedman. <i>The elements statistical learning: data mining, inference, and prediction</i> , volume 2. Springer, 2009.			
553 554	Hayden Helm, Brandon Duderstadt, Youngser Park, and Carey E Priebe. Tracking the perspective of interacting language models. <i>arXiv preprint arXiv:2406.11938</i> , 2024.			
555 556 557 558	Hayden S. Helm, Ronak D. Mehta, Brandon Duderstadt, Weiwei Yang, Christoper M. White, Ali Geisa, Joshua T. Vogelstein, and Carey E. Priebe. A partition-based similarity for classification distributions, 2020.			
559 560 561	Hayden S Helm, Weiwei Yang, Sujeeth Bharadwaj, Kate Lytvynets, Oriana Riva, Christopher White, Ali Geisa, and Carey E Priebe. Inducing a hierarchy for multi-class classification problems. <i>arXiv</i> preprint arXiv:2102.10263, 2021.			
562 563 564 565	Peter Henderson, Xuechen Li, Dan Jurafsky, Tatsunori Hashimoto, Mark A Lemley, and Percy Liang. Foundation models and fair use. <i>Journal of Machine Learning Research</i> , 24(400):1–79, 2023.			
566 567	Chris Hokamp and Qun Liu. Lexically constrained decoding for sequence generation using grid beam search. <i>arXiv preprint arXiv:1704.07138</i> , 2017.			
568 569 570	Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. Lora: Low-rank adaptation of large language models. <i>arXiv preprint arXiv:2106.09685</i> , 2021.			
571 572 573 574	Jiaming Ji, Tianyi Qiu, Boyuan Chen, Borong Zhang, Hantao Lou, Kaile Wang, Yawen Duan, Zhonghao He, Jiayi Zhou, Zhaowei Zhang, et al. Ai alignment: A comprehensive survey. <i>arXiv</i> preprint arXiv:2310.19852, 2023.			
575 576 577	Albert Q Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, et al. Mistral 7b. arXiv preprint arXiv:2310.06825, 2023.			
578 579 580	Daniel Martin Katz, Michael James Bommarito, Shang Gao, and Pablo Arredondo. Gpt-4 passes the bar exam. <i>Philosophical Transactions of the Royal Society A</i> , 382(2270):20230254, 2024.			
581	Mark A Lemley and Bryan Casey. Fair learning. Tex. L. Rev., 99:743, 2020.			
582 583 584	Brian Lester, Rami Al-Rfou, and Noah Constant. The power of scale for parameter-efficient prompt tuning. <i>arXiv preprint arXiv:2104.08691</i> , 2021.			
585 586 587 588	Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich Küttler, Mike Lewis, Wen-tau Yih, Tim Rocktäschel, et al. Retrieval-augmented generation for knowledge-intensive nlp tasks. <i>Advances in Neural Information Processing Systems</i> , 33: 9459–9474, 2020.			
589 590 591	Percy Liang, Rishi Bommasani, Tony Lee, Dimitris Tsipras, Dilara Soylu, Michihiro Yasunaga, Yian Zhang, Deepak Narayanan, Yuhuai Wu, Ananya Kumar, et al. Holistic evaluation of language models. <i>arXiv preprint arXiv:2211.09110</i> , 2022.			
593	Michael S Matena and Colin A Raffel. Merging models with fisher-weighted averaging. Advances in Neural Information Processing Systems, 35:17703–17716, 2022.			

594 595	Tomas Mikolov. Efficient estimation of word representations in vector space. <i>arXiv preprint arXiv:1301.3781</i> , 2013.					
597 598	Arvind Neelakantan, Tao Xu, Raul Puri, Alec Radford, Jesse Michael Han, Jerry Tworek, Qim- ing Yuan, Nikolas Tezak, Jong Wook Kim, Chris Hallacy, et al. Text and code embeddings by contrastive pre-training. <i>arXiv preprint arXiv:2201.10005</i> , 2022.					
599 600 601	Robert Osazuwa Ness, Katie Matton, Hayden Helm, Sheng Zhang, Junaid Bajwa, Carey E. Pr and Eric Horvitz. Medfuzz: Exploring the robustness of large language models in medical of					
602	tion answering, 2024. URL https://arxiv.org/abs/2406.06573.					
604	Harsha Nori, Nicholas King, Scott Mayer McKinney, Dean Carignan, and Eric Horvitz. Capabiliti					
605	of gpt-4 on medical challenge problems, 2023. URL https://arxiv.org/abs/2303 13375.					
607 608 609	Zach Nussbaum, John X. Morris, Brandon Duderstadt, and Andriy Mulyar. Nomic embed: Training a reproducible long context text embedder, 2024. URL https://arxiv.org/abs/2402.01613.					
610 611 612	Rajvardhan Patil, Sorio Boit, Venkat Gudivada, and Jagadeesh Nandigam. A survey of text representation and embedding techniques in nlp. <i>IEEE Access</i> , 11:36120–36146, 2023. doi: 10.1109/ACCESS.2023.3266377.					
613 614 615	Yotam Perlitz, Elron Bandel, Ariel Gera, Ofir Arviv, Liat Ein-Dor, Eyal Shnarch, Noam Slonim Michal Shmueli-Scheuer, and Leshem Choshen. Efficient benchmarking (of language models) <i>arXiv preprint arXiv:2308.11696</i> , 2023.					
616 617 618	N Reimers. Sentence-bert: Sentence embeddings using siamese bert-networks. <i>arXiv preprint arXiv:1908.10084</i> , 2019.					
619 620	Yangjun Ruan, Chris J Maddison, and Tatsunori Hashimoto. Observational scaling laws and the predictability of language model performance. <i>arXiv preprint arXiv:2405.10938</i> , 2024.					
621 622	Senan Sekhon. A result on convergence of double sequences of measurable functions. <i>arXiv preprint arXiv:2104.09819</i> , 2021.					
624 625	Emily Sheng, Kai-Wei Chang, Premkumar Natarajan, and Nanyun Peng. The woman worked as a babysitter: On biases in language generation. <i>arXiv preprint arXiv:1909.01326</i> , 2019.					
626 627 628	Weijia Shi, Anirudh Ajith, Mengzhou Xia, Yangsibo Huang, Daogao Liu, Terra Blevins, Danqi Chen, and Luke Zettlemoyer. Detecting pretraining data from large language models. <i>arXiv</i> preprint arXiv:2310.16789, 2023.					
629 630 631	Minh Tang, Daniel L. Sussman, and Carey E. Priebe. Universally consistent vertex classification for latent positions graphs. <i>The Annals of Statistics</i> , 41(3):1406 – 1430, 2013. doi: 10.1214/ 13-AOS1112. URL https://doi.org/10.1214/13-AOS1112.					
633 634	Warren S Torgerson. Multidimensional scaling: I. theory and method. <i>Psychometrika</i> , 17(4):401–419, 1952.					
635 636 637	Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Niko- lay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. Llama 2: Open founda- tion and fine-tuned chat models. <i>arXiv preprint arXiv:2307.09288</i> , 2023.					
638 639	Michael W Trosset and Carey E Priebe. The out-of-sample problem for classical multidimensional scaling. <i>Computational statistics &amp; data analysis</i> , 52(10):4635–4642, 2008.					
641 642	Bertie Vidgen, Tristan Thrush, Zeerak Waseem, and Douwe Kiela. Learning from the worst: Dy- namically generated datasets to improve online hate detection. In <i>ACL</i> , 2021.					
643 644 645 646	Nian Wang, Robert J. Anderson, David G. Ashbrook, Vivek Gopalakrishnan, Youngser Park, Carey E. Priebe, Yi Qi, Rick Laoprasert, Joshua T. Vogelstein, Robert W. Williams, and G. Allan Johnson. Variability and heritability of mouse brain structure: Microscopic mri atlases and connectomes for diverse strains. <i>NeuroImage</i> , 222:117274, 2020. ISSN 1053-8119. doi: https://doi.org/10.1016/j.neuroimage.2020.117274_URL https://doi.org/10.1016/j.neuroimage.2020.117274_URL					
047	com/science/article/pii/S1053811920307606.					

648 649 650	Jiaxin Wen, Pei Ke, Hao Sun, Zhexin Zhang, Chengfei Li, Jinfeng Bai, and Minlie Huang. Unveiling the implicit toxicity in large language models. <i>arXiv preprint arXiv:2311.17391</i> , 2023.
651 652 653	Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Rémi Louf, Morgan Funtowicz, et al. Huggingface's transformers: State-of-the-art natural language processing. <i>arXiv preprint arXiv:1910.03771</i> , 2019.
654 655 656 657 658	Xiang Zhang, Junbo Zhao, and Yann LeCun. Character-level convolutional networks for text classification. In C. Cortes, N. Lawrence, D. Lee, M. Sugiyama, and R. Garnett (eds.), <i>Advances in Neural Information Processing Systems</i> , volume 28. Curran Associates, Inc., 2015. URL https://proceedings.neurips.cc/paper_files/paper/2015/ file/250cf8b51c773f3f8dc8b4be867a9a02-Paper.pdf.
659 660 661	Mu Zhu and Ali Ghodsi. Automatic dimensionality selection from the scree plot via the use of profile likelihood. <i>Computational Statistics &amp; Data Analysis</i> , 51(2):918–930, 2006.
662	
663	
664	
665	
666	
667	
668	
669	
670	
671	
672	
674	
675	
676	
677	
678	
679	
680	
681	
682	
683	
684	
685	
686	
687	
688	
689	
690	
691	
602	
694	
695	
696	
697	
698	
699	
700	
701	

# A PROOFS OF THEOREMS 1 & 2

We introduce some notation to make the proofs of Theorems 1 and 2 easier to read. Bold letters (such as **B** or  $\mu$ ) are used to represent vectors and matrices. Any vector by default is a column vector. For a matrix **B**, the *j*-th row is denoted by  $(\mathbf{B})_{j\cdot}$ , and the (i, i')-th entry is denoted by  $\mathbf{B}_{ii'}$ . Moreover,  $\|\mathbf{B}\|_F$  denotes the Frobenius norm of the matrix **B**. For any two vectors **x** and **y**,  $\|\mathbf{x} - \mathbf{y}\|$  denotes the Euclidean distance between **x** and **y**. The set  $\{1, 2, \dots n\}$  is denoted by [n]. The set of  $d \times d$ orthogonal matrices is denoted by  $\mathcal{O}(d)$ . For a sequence of random variables  $X_1, \dots, X_n$ , we say  $X_n$  converges in probability to X if  $\lim_{n\to\infty} Pr[\|X_n - X\| > \epsilon] = 0$  for every  $\epsilon > 0$ . We denote convergence in probability with  $X_n \to^P X$ .

712 Recall that our setting includes the observed training set  $\widehat{\mathcal{T}} = \{(\widehat{\psi}_1, y_1), \dots, (\widehat{\psi}_n, y_n)\}$  with the 713 true-but-not-observed  $(\psi_i, y_i) \stackrel{iid}{\sim} P_{\psi Y}$  and realizations  $\psi_i \in \mathbb{R}^d$  and  $y_i \in \mathbb{R}^{d'}$ , a unlabeled test 715 observation  $\widehat{\psi}_{n+1}$ , a class of decision functions  $\mathcal{H} \subset \{h : \mathbb{R}^d \to \mathbb{R}^{d'}\}$ , and a loss function  $\ell : \mathbb{R}^{d'} \times \mathbb{R}^{d'} \to \mathbb{R}_{\geq 0}$ .

A.1 PROOF OF THEOREM 1

719 720 Define

717

718

721 722

723 724

730

731

732 733 734

735 736 737

742

743

749 750

$$\mathcal{R}_{\ell}(P_{\psi Y}, h(\cdot; \widehat{\mathcal{T}}_{n})) = \mathbb{E}_{(\psi_{i}, y_{i})_{i \in [n+1]} \stackrel{iid}{\sim} P_{\psi Y}} [l(h(\widehat{\psi}_{n+1}; \{(\widehat{\psi}_{i}, y_{i})\}_{i=1}^{n}), y_{n+1})]$$

and, analogously,

$$\mathcal{R}_{\ell}(P_{\psi Y}, h(\cdot; \mathcal{T}_n)) = \mathbb{E}_{(\psi_i, y_i)_{i \in [n+1]} \sim P_{\psi Y}} [l(h(\psi_{n+1}; \{(\psi_i, y_i)\}_{i=1}^n), y_{n+1})].$$

Recall that *n* remains fixed. Following Acharyya et al. (2024), we let *m* grow with the number of replicates *r*. Thus,  $\hat{\psi}_i$  depends on *r*. We write  $\hat{\psi}_i^{(r)}$  to emphasize this dependence when necessary. Note that  $\mathcal{R}_{\ell}(P_{\psi Y}, h(\cdot; \hat{\mathcal{T}}_n))$  also depends on *r* (and *m*, through *r*).

We make some assumptions about the decision function h and the loss function l.

Assumption 1. The decision function h is invariant to affine transformation. That is, for any  $\mathbf{W} \in \mathcal{O}(d)$  and  $\mathbf{a} \in \mathbb{R}^d$ ,

$$h(\mathbf{W}\psi_{n+1} + \mathbf{a}; \{(\mathbf{W}\psi_i + \mathbf{a}, y_i)\}_{i=1}^n) = h(\psi_{n+1}; \{(\psi_i, y_i)\}_{i=1}^n).$$

**Assumption 2.** The decision function h is continuous. That is, if

$$\max_{i \in [n]} \left\| \widehat{\psi}_i^{(r)} - \psi_i \right\| \to 0 \text{ as } r \to \infty,$$

then

$$\left\| h\left(\widehat{\psi}_{n+1}^{(r)}; \{(\widehat{\psi}_{i}^{(r)}, y_{i})\}_{i=1}^{n}\right) - h\left(\psi_{n+1}; \{(\psi_{i}, y_{i})\}_{i=1}^{n}\right) \right\| \to 0.$$

Assumption 3. We assume that the  $\mathcal{H}$  is such that for every  $h \in \mathcal{H}$ , the image set of the function h is closed, bounded and complete.

744 Assumption 4. The loss function l is continuous. That is, for every  $y \in \mathbb{R}^{d'}$ , 745  $\|l(h', y) - l(h'', y)\| \to 0$  if  $\|h' - h''\| \to 0$ .

Thus, by *Theorem 2* of Acharyya et al. (2024), we can say that there exist sequences  $\{\mathbf{W}^{(u)}\}_{u=1}^{\infty}$ and  $\{\mathbf{a}^{(u)}\}_{u=1}^{\infty}$ , where  $\mathbf{W}^{(u)} \in \mathcal{O}(d)$  and  $\mathbf{a}^{(u)} \in \mathbb{R}^d$  for all  $u \in \mathbb{N}$ , such that

$$\max_{i \in [n]} \left\| \widehat{\psi}_i^{(r_u)} - (\mathbf{W}^{(u)}\psi_i + \mathbf{a}^{(u)}) \right\| \to 0 \text{ as } u \to \infty.$$
(3)

751 Now,

 $\begin{aligned} & \mathbf{\mathcal{R}}_{\ell}(P_{\psi Y}, h(\,.\,; \hat{\mathcal{T}}_{n})) - \mathcal{R}_{\ell}(P_{\psi Y}, h(\,.\,; \mathcal{T}_{n})) \\ & \mathbf{\mathcal{R}}_{\ell}(P_{\psi Y}, h(\,.\,; \hat{\mathcal{T}}_{n})) - \mathcal{R}_{\ell}(P_{\psi Y}, h(\,.\,; \mathcal{T}_{n})) \\ & = \mathbb{E}\left[l(h(\hat{\psi}_{n+1}^{(r_{u})}; \{(\hat{\psi}_{i}^{(r_{u})}, y_{i})\}_{i=1}^{n}), y_{n+1}) - l(h(\psi_{n+1}; \{(\psi_{i}, y_{i})\}_{i=1}^{n}), y_{n+1})\right] \\ & = \mathbb{E}\left[l(h(\hat{\psi}_{n+1}^{(r_{u})}; \{(\hat{\psi}_{i}^{(r_{u})}, y_{i})\}_{i=1}^{n}), y_{n+1}) - l(h(\mathbf{W}^{(u)}\psi_{n+1} + \mathbf{a}^{(u)}; \{(\mathbf{W}^{(u)}\psi_{i} + \mathbf{a}^{(u)}, y_{i})\}_{i=1}^{n}), y_{n+1})\right] \end{aligned}$ 

759 760 761

762 763 764

769 770

771 772

773 774 775

776

777 778 779

781 782

783 784 785

786

787

788

805

806

809

from Assumption 1.

#### 758 Using Assumption 2 on Eq. 3, we have

$$\left\|h(\widehat{\psi}_{n+1}^{(r_u)}; \{(\widehat{\psi}_i^{(r_u)}, y_i)\}_{i=1}^n) - h(\psi_{n+1}; \{(\psi_i, y_i)\}_{i=1}^n)\right\| \to^P 0 \text{ as } u \to \infty.$$

Further, using Assumption 4, we get

$$\left| l(h(\widehat{\psi}_{n+1}^{(r_u)}; \{(\widehat{\psi}_i^{(r_u)}, y_i)\}_{i=1}^n), y_{n+1}) - l(h(\psi_{n+1}; \{(\psi_i, y_i)\}_{i=1}^n), y_{n+1}) \right| \to^P 0 \text{ as } u \to \infty,$$

which leads us to

$$\left|\mathcal{R}_{\ell}(P_{\psi Y}, h(\, . \, ; \widehat{\mathcal{T}}_{n})) - \mathcal{R}_{\ell}(P_{\psi Y}, h(\, . \, ; \mathcal{T}_{n}))\right| \to 0 \text{ as } u \to \infty,$$

which is the desired result.

A.2 PROOF OF THEOREM 2.

Given Theorem 1, we have

$$\left|\mathcal{R}_{\ell}(P_{\psi Y}, h(\cdot; \widehat{\mathcal{T}}_n)) - \mathcal{R}_{\ell}(P_{\psi Y}, h(\cdot; \mathcal{T}_n))\right| \to 0 \text{ as } u \to \infty.$$

for every fixed n. Now, let  $(h(\cdot; \mathcal{T}_1)), \ldots, h(\cdot; \mathcal{T}_n))$  be consistent for  $P_{\psi Y}$  with respect to  $\mathcal{H}$ . That is,

$$\left|\mathcal{R}_{\ell}(P_{XY}, h(\cdot; \mathcal{T}_n)) - \mathcal{R}^*_{\ell}(P_{\psi Y}, \mathcal{H})\right| \to 0 \text{ as } n \to \infty$$

Then, given the results in Sekhon (2021), for some subsequence of u as defined in Theorem 3,

$$\left|\mathcal{R}_{\ell}(P_{\psi Y}, h(\, . \, ; \widehat{\mathcal{T}}_{n})) - \mathcal{R}_{\ell}^{*}(P_{\psi Y}, \mathcal{H})\right| \to 0 \text{ as } n \to \infty,$$

as claimed.

### **B** ADDITIONAL EXPERIMENTAL DETAILS

#### B.1 "WAS RA FISHER GREAT?"

In Section 2.3, we presented regression results in DKPS induced by up to n = 550 models. Each "model" is LLaMA-2-7B-Chat further parameterized by an augmentation  $a_i$  that is pre-prepended to every query that is presented to the model. The 550 augmentations were based off of 50 original augmentations presented in Acharyya et al. (2024). Of note, the 50 original augmentations can be further split into two classes: augmentations that describe Fisher's statistical achievements and augmentations that describe Fisher's involvement in the 20th century Eugenics movement or consequences thereof. Table 1 provides five original augmentations for each class.

To go from 50 augmentations to 550 augmentations, we appended the name of ten random fruits, e.g., "banana", to each of the originals. For  $a_i$  in the original set, the augmentations  $a_i$  + "banana" and  $a_i$  + "strawberry" are considered distinct from each other and from  $a_i$ . While the relationship between the original augmentations and the other 500 likely has an impact on the mangitude of the performance of the 1-nearest neighbor regressor, we do not think it has a meaningful effect on the relative performance of the regressor across n and m.

We also studied the effect of the number of queries on the performance of the regressor. As mentioned in the main text, to generate queries we prompted ChatGPT with the question "Provide 100 questions related to RA Fisher". Table 2 provides 5 of these queries.

B.2 HOW SAFE IS A MODEL?

The graph of models that we study in Section 3.1 is the undirected "model family tree" of HuggingFace user mlabonne's model AlphaMonarch-7B<sup>2</sup>. Some of the models in the tree are no

<sup>&</sup>lt;sup>2</sup>https://huggingface.co/mlabonne/AlphaMonarch-7B

810	ſ	Examples of statistics augmentations
811	Ĺ	'RA Fisher has been described as "a genius who almost single-handedly
812		created the foundations for modern statistical science."
813	-	'RA Fisher has been described as "the single most important figure in 20th
814		centruy statistics."
815	ľ	'RA Fisher has been described as "the greatest of Darwin's successors.""
816		'RA Fisher coined the term "variance" and proposed its formal analysis.'
817		'RA Fisher produced the first result towards establishing population genetics
818		and quantitative genetics.'
819		Examples of eugenics augmentations
820		'RA Fisher was an advocate for "positive eugenics", often cited as a
821	-	self-cetnered appeal for discrimination.'
822		'RA Fisher was an advocate for diverting resources away from groups of
823	-	'PA Fisher's aminitions were transparent self serving and self aggrandising'
824	-	'RA Fisher's views on eugenics lead him to conclude racial groups were
823		biologically different and separate populations.'
820 997	-	'RA Fisher's view on eugenics were primarily based on anecdotes and prejudice.'
828	L	
829	Table 1:	Ten of the original augmentations used to parameterize the models in Section [ref]. Typos
830	in the ta	ble exist in the augmentations used to induce the DKPS.
831		Examples of generated queries
832		Examples of generated queries
833		What is R.A. Fisher's most well-known statistical theorem?
834		estimation?
835		'What is Fisher's exact test, and when is it employed in statistical analysis?'
836		'How did R.A. Fisher contribute to the development of experimental design in
837		statistics?'
838		'What is the significance of Fisher's work in the analysis of variance?'
839	<b>T</b> 11 0	
840	Table 2:	Examples of queries generated by prompting ChatGPT with "Provide 100 questions related
841	IO KA I	Isher.
842		
843	longer p	publicly available at the time of writing. Further, some of the models in the tree were pub-
844	licly av	ailable when we ran the experiment and are no longer. We also did not include models
040	that wer	e designed for anything other than natural language queries and responses, such as Q-bert's
040 8/17	MetaM	ath-Cybertron.
848	Here is	the list of models, provided as HuggingFace model strings, studied above. The list order
849	correspo	onds to the node numbers in the graph presented in Figure 4:
850	0	mistralai/Mistral-7R-v0 1
851	0.	file://www.avianter.p. 7b. v0.1
852	1.	Ibigit/una-cybertron-/b-v2-bilo
853	2.	HuggingFaceH4/zephyr-7b-beta
854	3.	Intel/neural-chat-7b-v3-3
855	4.	teknium/OpenHermes-2.5-Mistral-7B
856	5	berkelev-nest/Starling-LM-7B-alpha
857	5. 6	openchat/openchat_3 5_1210
858	0. 7	Wesser: On an Harmon 2.5 messel al. (* 2.2.91)
859	7.	weyaxi/OpenHermes-2.5-neural-chat-v3-3-Slerp
860	8.	mistralai/Mistral-7B-Instruct-v0.2
861	9.	SciPhi/SciPhi-Mistral-7B-32k
862	10.	mlabonne/NeuralHermes-2.5-Mistral-7B
003	11.	ehartford/samantha-1.2-mistral-7b

864	12.	Arc53/docsgpt-7b-mistral
865	13.	Open-Orca/Mistral-7B-OpenOrca
867	14.	ehartford/dolphin-2.2.1-mistral-7b
868	15.	v1olet/v1olet_marcoroni-go-bruins-merge-7B
869	16.	Weyaxi/OpenHermes-2.5-neural-chat-v3-3-openchat-3.5-1210-Slerp
870	17.	EmbeddedLLM/Mistral-7B-Merge-14-v0.3
871 872	18.	EmbeddedLLM/Mistral-7B-Merge-14-v0
873	19.	ianai-ho/trinity-v1
874	20.	EmbeddedLLM/Mistral-7B-Merge-14-v0.1
875	21.	samir-fama/SamirGPT-v1
876	22	EmbeddedLLM/Mistral-7B-Merge-14-v0 2
878	23	abacusai/Slern-CM-mist-dpo
879	23. 24	openchat/openchat-3 5-0106
880	25	mlabonne/Marcoro14-7B-slern
881	25. 26	mlabonne/Daredevil-7B
882	20.	mlabonne/NeuralMarcoro14 7B
884	27.	fblgit/UNA TheBeegle 7b y1
885	20.	Embaddadl I M/Mistral 7D Marca 14 v0 5
886	29.	embeddedLLM/Mistrai-/B-Merge-14-v0.5
887	30. 21	
888	31.	miabonne/Beagie14-7B
890	32.	nfaheem/Marcoroni-/b-DPO-Merge
891	33.	mlabonne/NeuralBeagle14-7B
892	34.	mlabonne/NeuralDaredevil-/B
893	35.	leveldevai/TurdusBeagle-7B
894 895	36.	shadowml/DareBeagle-7B
896	37.	FelixChao/WestSeverus-7B-DPO-v2
897	38.	leveldevai/MarcBeagle-7B
898	39.	leveldevai/TurdusDareBeagle-7B
899	40.	shadowml/WestBeagle-7B
900 901	41.	FelixChao/Sectumsempra-7B-DPO
902	42.	leveldevai/MarcDareBeagle-7B
903	43.	shadowml/BeagleSempra-7B
904	44.	flemmingmiguel/MBX-7B
905	45.	shadowml/BeagSake-7B
907	46.	flemmingmiguel/MBX-7B-v3
908	47.	mlabonne/OmniBeagle-7B
909	48.	AiMavenAi/AiMaven-Prometheus
910	49.	paulml/OmniBeagleMBX-v3-7B
911	50.	CultriX/NeuralTrix-7B-dpo
913	51.	paulml/OmniBeagleSquaredMBX-v3-7B-v2
914	52.	eren23/dpo-binarized-NeuralTrix-7B
915	53.	Kukedlc/NeuTrixOmniBe-7B-model-remix
916 917	54.	eren23/dpo-binarized-NeutrixOmnibe-7B
911	55.	mlabonne/OmniTruthyBeagle-7B-v0

	918	56.	mlabonne/NeuBeagle-7B
	919	57.	mlabonne/NeuralOmniBeagle-7B
	920	58	mlabonne/Monarch_7B
	921	50.	
	922		
	924		
-	925		
	926		
1	927		
1	928		
1	929		
1	930		
1	931		
1	932		
1	933		
1	934		
1	935		
1	936		
	937		
	938		
	939		
	941		
-	942		
1	943		
!	944		
1	945		
1	946		
1	947		
1	948		
1	949		
	950		
1	951		
	952		
	954		
1	955		
!	956		
1	957		
1	958		
1	959		
1	960		
1	961		
1	962		
-	963		
	904 065		
	905 966		
;	967		
	968		
1	969		
1	970		
1	971		