

Honest-but-Curious Nets: Sensitive Attributes of Private Inputs Can Be Secretly Coded into the Classifiers’ Outputs

Mohammad Malekzadeh
Imperial College London, UK
m.malekzadeh@imperial.ac.uk

Anastasia Borovykh
Imperial College London, UK
a.borovykh@imperial.ac.uk

Deniz Gündüz
Imperial College London, UK
d.gunduz@imperial.ac.uk

ABSTRACT

It is known that deep neural networks, trained for the classification of non-sensitive *target* attributes, can reveal *sensitive* attributes of their input data through internal representations extracted by the classifier. We take a step forward and show that deep classifiers can be trained to secretly encode a sensitive attribute of their input data into the classifier’s outputs for the target attribute, at inference time. Our proposed attack works even if users have a full white-box view of the classifier, can keep all internal representations hidden, and only release the classifier’s estimations for the target attribute. We introduce an information-theoretical formulation for such attacks and present efficient empirical implementations for training *honest-but-curious* (HBC) classifiers: *classifiers that can be accurate in predicting their target attribute, but can also exploit their outputs to secretly encode a sensitive attribute*. Our work highlights a vulnerability that can be exploited by malicious machine learning service providers to attack their user’s privacy in several seemingly safe scenarios; such as encrypted inferences, computations at the edge, or private knowledge distillation. Experimental results on several attributes in two face-image datasets show that a semi-trusted server can train classifiers that are not only perfectly *honest* but also accurately *curious*. We conclude by showing the difficulties in distinguishing between standard and HBC classifiers, discussing challenges in defending against this vulnerability of deep classifiers, and enumerating related open directions for future studies.

CCS CONCEPTS

• **Security and privacy** → **Privacy protections; Information-theoretic techniques; Social aspects of security and privacy.**

KEYWORDS

privacy in machine learning; attacks on privacy; data privacy

1 INTRODUCTION

Machine learning (ML) classifiers, trained on a set of labeled data, aim to facilitate the estimation of a *target* label (*i.e.*, *attribute*) for new data at inference time; from smile detection for photography [82] to automated detection of a disease on medical data [12, 19]. However, in addition to the target attribute that the classifier is trained for, data might also contain some other *sensitive* attributes. For example, there are several attributes that can be inferred from a face image; such as gender, age, race, emotion, hairstyle, and more [35, 86].

Will be appeared in ACM Conference on Computer and Communications Security (CCS ’21), November 15–19, 2021.

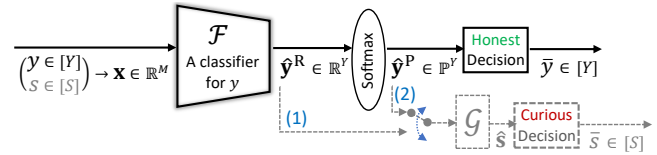


Figure 1: Let users’ data x (e.g., face image) contain a *target* attribute y (e.g., Age) and a *sensitive* attribute s (e.g., Race). A classifier \mathcal{F} is provided by a server to estimates y , and users only release the classifier’s outputs \hat{y} . We show that \mathcal{F} can be trained such that \hat{y} is not only accurate for y (*honesty*), but can also be used by a secret attack \mathcal{G} to infer s (*curiosity*). We present efficient attacks in two scenarios where either (1) the raw outputs \hat{y}^R , or (2) the soft outputs \hat{y}^P are released.

Since ML classifiers, particularly deep neural networks (DNNs), are becoming increasingly popular, either as cloud-based services or as part of apps on our personal devices, it is important to be aware of the type of sensitive attributes that we might reveal through using these classifiers; especially when a classifier is supposed to only estimate a specified attribute. For instance, while clinical experts can barely identify the race of patients from their medical images, DNNs show considerable performance in detecting race from chest X-rays and CT scans [3].

In two-party computations, a legitimate party that does not deviate from its specified protocol but attempts to infer as much sensitive information as possible from the received data is called a *honest-but-curious* (HBC) party [9, 17, 30, 54]. Following convention, if a classifier’s outputs not only allow to estimate the target attribute, but also reveal information about other attributes (particularly those uncorrelated to the target one) we call it an *HBC classifier*. In this paper, we show how a semi-trusted server can train an HBC classifier such that the outputs of the classifier are not only useful for inferring the target attribute, but can also secretly carry information about a sensitive attribute of the user’s data that is unrelated to the target attribute.

Figure 1 shows an overview of the problem. We consider a server that provides its users access to a classifier trained for a known target attribute y . We put no restriction on the input data x or the users’ access to the classifier: users can have control over the classifier and get a white-box view to it (e.g., when the classifier is deployed on the users’ devices), or users can perform secure computation on their data if they have a black-box view (e.g., when the classifier is hosted in the cloud). Our only assumption is that both the user and the server can observe the outputs of classifier \hat{y} . For a Y -class classifier, \hat{y} is a real-valued vector of size Y containing either (1) raw

scores $\hat{y}^R \in \mathbb{R}^Y$, or (2) *soft* scores $\hat{y}^P = \text{softmax}(\hat{y}^R) \in \mathbb{P}^Y$, where the i -th soft score is $\hat{y}_i^P = e^{\hat{y}_i^R} / \sum_{j=0}^{Y-1} e^{\hat{y}_j^R}$. The server usually uses some threshold functions to decide \bar{y} as the final predicted class. In classification tasks, an estimated probability distribution over possible classes is more useful than just receiving the most probable class; as it allows the aggregation of outputs provided by multiple ML services to enhance the ultimate decision. Moreover, collecting outputs can help a server to monitor and enhance its decisions and the provided services.

To protect the *user's* privacy, it is usually proposed to hide the data, as well as all the intermediate computations on the data, and only release the classifier's outputs; either using secure two-party computation via cryptography [2, 4, 16] or by restricting the computations to edge devices [45, 69]. Although these encrypted or edge solutions hide the input data and the internal features extracted by the classifier, the outputs are usually released to the service provider, because the estimation of a target attribute might not seem sensitive to the user's privacy and it might be needed for further services offered to the user. For instance, an insurance company raised huge ethical concerns, when it announced that its ML model extracts "non-verbal cues" from videos of users' faces to identify fraud [33]. We show that even if the video is processed locally, or in an encrypted manner, and only a single real-valued output $\hat{y} \in [0, 1]$ is released to the insurance company, as the probability of fraud, this single output can still be designed to reveal another sensitive attribute about the user. In our experiments, we show that a $\hat{y} \in [0, 1]$ produced by a smile-detection classifier, as the probability of "smiling", can be used to secretly infer whether that person is "white" or not¹.

We first show that in a *black-box* view, where an arbitrary architecture can be used for the classifier, the server can obtain the best achievable trade-off between *honesty* (i.e., the classification accuracy for the target attribute) and *curiosity* (i.e., the classification accuracy for the sensitive attribute). Specifically, we build a controlled synthetic dataset and show how to create such an HBC classifier via a weighted mixture of two separately trained classifiers, one for the target attribute and another for the sensitive attribute (Section 3). Then, we focus on the more challenging, *white-box* view where the server might have some constraints on the chosen model, e.g., the restriction to not being suspicious, or that the classifier must be one of the known off-the-shelf models (Section 4). To this end, we formulate the problem of training an HBC classifier in a general information-theoretical framework, via the *information bottleneck* principle [71], and show the existence of a general attack for encoding a desired sensitive data into the output of a classifier. We propose two practical methods that can be used by a server for building an HBC classifier, one via the *regularization* of classifier's loss function, and another via training of a *parameterized* model.

Extensive experiments, using typical DNNs for several tasks with different attributes defined on two real-world datasets [35, 86], show that HBC classifiers can mostly achieve honesty very close to standard classifiers, while also being very successful in their curiosity (Section 5). We, theoretically and empirically, show that the entropy of an HBC classifier's outputs usually tends to be higher

than the entropy of a standard classifier's outputs. Moreover, we explain how a server can improve the honesty of the classifier by trading some curiosity via adding an entropy minimization component to parameterized attacks, which in particular, can make HBC classifiers less suspicious against proactive defenses [28].

Previous works propose several types of attacks to ML models [10, 44], mostly to DNNs [34], including property inference [43], membership inference [59, 63], model inversion [15], model extraction [25, 73], adversarial examples [67], or model poisoning [7, 26]. But these attacks mostly concern the privacy of the training dataset and, in all these attacks, the ML model is the trusted party, while users are assumed untrusted. Our work, from a different point of view, discusses a new threat model, where (the owner of) the ML model is semi-trusted and might attack the privacy of its users at inference time. The closest related work is the "overlearning" concept in [66], where it is shown that internal representations extracted by DNN layers can reveal sensitive attributes of the input data that might not even be correlated to the target attribute. The assumption of [66] is that an adversary observes a subset of internal representations, while we assume all internal representations to be hidden and an adversary has access only to the outputs. Notably, we show that when users only release the classifier's outputs, overlearning is not a major concern as standard classifiers do not reveal significant information about a sensitive attribute through their outputs, whereas an HBC classifier can secretly, and almost perfectly, reveal a sensitive attribute just via classifier's outputs.

Contributions. In summary, this paper proposes the following contributions to advance privacy protection in using ML services.

- (1) We show that ML services can attack their user's privacy even in a highly restricted setting where they can only get access to the results of an agreed target computation on their users' data.
- (2) We formulate such an attack in a general information-theoretical formulation and show the efficiency of our attack via several empirical results. Mainly, we show how HBC classifiers can encode a sensitive attribute of their private input into the classifier's output, by exploiting the output's entropy as a side-channel. Therefore, HBC classifiers tend to produce higher-entropy outputs than standard classifiers. However, we also show that the output's entropy can be efficiently reduced by trading a small amount of curiosity of the classifier, thus making it even harder to distinguish standard and HBC classifiers.
- (3) We show an important threat of this vulnerability in a recent approach where knowledge distillation [22] is used to train a student classifier on private unlabeled data via a teacher classifier that is already trained on a set of labeled data, and show that an HBC teacher can transfer its curiosity capability to the student classifiers.
- (4) We support our findings via several experimental results on two real-world datasets with different characteristics, as well as additional analytical results for the setting of training convex classifiers.

Code and instructions for reproducing the reported results are available at <https://github.com/mmalekzadeh/honest-but-curious-nets>.

¹See Appendix E for more motivational examples of HBC classifiers that can be trained for other types of users' data, such as text, motion sensors, and audio.

2 PROBLEM FORMULATION

Notation. We use lower-case *italic*, e.g., x , for scalar variables; upper-case *italic*, e.g., X , for scalar constants; lower-case bold, e.g., \mathbf{x} , for vectors; upper-case blackboard, e.g., \mathbb{X} , for sets; calligraphic font, e.g., \mathcal{X} , for functions; subscripts, e.g., w_1 , for indexing a vector; superscripts, e.g., w^1 , for distinguishing different instances; \mathbb{R}^X for real-valued vectors of dimension X ; and \mathbb{P}^X for a probability simplex of dimension $X - 1$ (that denotes the space of all probability distributions on a X -value random variable). We have $\{0, \dots, X - 1\} \equiv [X]$, and $\lfloor \cdot \rfloor$ shows rounding to the nearest integer. Logarithms are natural unless written explicitly otherwise. The standard logistic is $\sigma(\mathbf{x}) = 1/(1 + \exp(-\mathbf{x}))$. Given random variables $a \in [X]$, $b \in [X]$, and $c \in [Z]$, the entropy of a is $H(a) = -\sum_{i=0}^{X-1} \Pr(a = i) \log \Pr(a = i)$, the cross entropy of b relative to a is $H_a(b) = -\sum_{i=0}^{X-1} \Pr(a = i) \log \Pr(b = i)$, the conditional entropy of a given c is $H(a|c) = -\sum_{i=0}^{X-1} \sum_{j=0}^{Z-1} \Pr(a = i, c = j) \log (\Pr(a = i, c = j)/\Pr(c = j))$, and the mutual information (MI) between a and c is $I(a; c) = H(a) - H(a|c)$ [38]. $\mathbb{I}_{(C)}$ shows the indicator function that outputs 1 if condition C holds, and 0 otherwise.

Definitions. Let a user own data $\mathbf{x} \in \mathbb{R}^M$ sampled from an unknown data distribution \mathcal{D} . Let \mathbf{x} be informative about at least two latent categorical variables (*attributes*): $y \in [Y]$ as the *target* attribute, and $s \in [S]$ as the *sensitive* attribute (see Figure 1). Let a server own a *classifier* \mathcal{F} that takes \mathbf{x} and outputs: $\hat{\mathbf{y}} = \mathcal{F}(\mathbf{x}) = [\hat{y}_0, \hat{y}_1, \dots, \hat{y}_{Y-1}]$, where \hat{y}_i estimates $\Pr(y = i|\mathbf{x})$. Let \bar{y} denote the predicted value for y that is decided from $\hat{\mathbf{y}}$; e.g., based on a *threshold* in binary classification or *argmax* function in multi-class classification. We assume that \mathbf{x} , and all intermediate computations of \mathcal{F} , are hidden and the user only releases $\hat{\mathbf{y}}$. Let $\hat{s} = \mathcal{G}(\hat{\mathbf{y}})$ be the *attack* on the sensitive attribute s that only the server knows about. Let \bar{s} denote the predicted value for s that is decided based on \hat{s} . In sum, the following Markov chain holds: $(y, s) \rightarrow \mathbf{x} \rightarrow \hat{\mathbf{y}} \rightarrow \hat{s}$.

Throughout this paper, we use the following terminology:

1. Honesty. Given a test dataset $\mathbb{D}^{test} \sim \mathcal{D}$, we define \mathcal{F} as a δ^y -*honest* classifier if

$$\Pr_{(\mathbf{x}, y) \sim \mathbb{D}^{test}, \bar{y} \leftarrow \mathcal{F}(\mathbf{x})} [\bar{y} = y] \geq \delta^y,$$

where $\delta^y \in [0, 1]$ is known as the classifier’s test accuracy, and we call it the *honesty* of \mathcal{F} in predicting the *target* attribute.

2. Curiosity. Given a test dataset $\mathbb{D}^{test} \sim \mathcal{D}$ and an attack \mathcal{G} , we define \mathcal{F} as a δ^s -*curious* classifier if

$$\Pr_{(\mathbf{x}, s) \sim \mathbb{D}^{test}, \bar{s} \leftarrow \mathcal{G}(\mathcal{F}(\mathbf{x}))} [\bar{s} = s] \geq \delta^s,$$

where $\delta^s \in [0, 1]$ is the attack’s success rate on the test set, and we call it the *curiosity* of \mathcal{F} in predicting the *sensitive* attribute.

3. Honest-but-Curious (HBC). We define \mathcal{F} as a (δ^y, δ^s) -HBC classifier if it is both δ^y -*honest* and δ^s -*curious* on the same \mathbb{D}^{test} .

4. Standard Classifier. A classifier \mathcal{F} that is trained only for achieving the best honesty, without any intended curiosity.

5. Black- vs. White-Box. We consider the users’ perspective to the classifier \mathcal{F} at *inference* time. In a *black-box* view, a user observes only the classifier’s outputs and not the classifier’s architecture and parameters. In a *white-box* view, a user also has full access to the classifier’s architecture, parameters, and intermediate computations.

6. Threat Model. The semi-trusted server chooses the algorithm and dataset ($\mathbb{D}^{train} \sim \mathcal{D}$) for training \mathcal{F} . In a black-box view, the server has the additional power to choose the architecture of \mathcal{F} (unlike the white-box view). At inference time, a user (who does not necessarily participate in the training dataset) runs the trained classifier on her private data once, and only reveals the classifier’s outputs $\hat{\mathbf{y}} = \mathcal{F}(\mathbf{x})$ to the server. We assume no other information is provided to the server at inference time.

Our Objective. We show how a server can train an HBC classifier to establish efficient *honesty-curiosity* trade-offs over achievable (δ^y, δ^s) pairs, and analyze the privacy risks, behavior, and characteristics of HBC classifiers compared to standard ones.

3 BLACK-BOX VIEW: A MIXTURE MODEL

To build a better intuition, we first discuss the *black-box* view where the server can choose any arbitrary architecture. and we show the existence of an efficient attack for every task with $S \leq Y$.

3.1 A Convex Classifier

We start with a simple logistic regression classifier. Let us consider the synthetic data distribution depicted in Figure 2, where each sample $\mathbf{x} \in \mathbb{R}^2$ has two attributes $y \in \{0, 1\}$ and $s \in \{0, 1\}$. If s is correlated with the y , the output of any classifier always reveals some sensitive information (which we explore it in real-world datasets in Section 5). The less s is correlated with the y , the more difficult it should be for the server to build an HBC classifier. Thus, the data distribution in Figure 2 is built such that y and s are independent, and for each attribute, there is an optimal linear classifier.

It is clear that for this dataset, we can find an optimal logistic regression classifier $\hat{y} = \sigma(w_1 x_1 + w_2 x_2 + b)$, with parameters $[w_1, w_2, b]$, that simulates the decision boundary of y . Such a classifier is δ^y -*honest* with $\delta^y = 1$, and at the same time it is δ^s -*curious* with $\delta^s = 0.5$; that means the classifier is honest and does not leak any sensitive information. The main point is that any effort for making a curious linear classifier with $\delta^s > 0.5$ will hurt the honesty by forcing $\delta^y < 1$. On this dataset, it can be shown that for any logistic regression classifier we have $\delta^y = 1.5 - \delta^s$. For example, the optimal linear classifier for attribute s cannot have a better performance than a random guess on attribute y .

In Appendix A, we show how a logistic regression classifier can become HBC with a convex loss function, and analyze the behavior of such a classifier in detail. Specifically, we show that the trade-off for a classifier with *limited* capacity (e.g., logistic regression) is that: if alongside the target attribute, we also optimize for the sensitive attribute, we will only ever converge to a neighborhood of the optimum for the target attribute. We show that the size of the neighborhood is getting larger by the weight (i.e., importance) we give to curiosity. While the analysis in Appendix A holds for a convex setting and are simplistic in nature, it provides intuitions into the idea that when the attributes y and s are somehow *correlated*, the output \hat{y} can better encode both tasks, but when we have *independent* attributes, we need classifiers with more capacity to cover the payoff for not converging to the optimal point of target attribute.

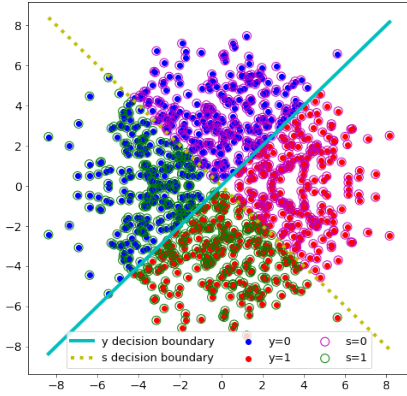


Figure 2: A two-attribute two-class data distribution on \mathbb{R}^2 , where attributes y and s are independent, and classes in each attribute are linearly separable.

3.2 A Mixture of Two Classifiers

Figure 3 shows how two logistic regression classifiers, each trained separately for a corresponding attribute, can be combined such that the final output \hat{y} is a mixture of the predicted values for the target attribute y , and the sensitive attribute s . There are two ways to combine $z^y \in [0, 1]$ and $z^s \in [0, 1]$. Considering multipliers $\beta^y \in [0, 1]$ and $\beta^s = 1 - \beta^y$, one option is the *normal* mixture, where $\hat{y} = \beta^y z^y + \beta^s z^s$, and another one is the *hard* mixture, where $\hat{y} = \beta^y \lfloor z^y \rfloor + \beta^s \lfloor z^s \rfloor$. Since two classifiers are each optimal for the dataset in Figure 2, \hat{y} in hard mixture can only take four values:

$$\hat{y} = \begin{cases} 0 & \text{if } y = 0 \text{ \& } s = 0 \\ \beta^s & \text{if } y = 0 \text{ \& } s = 1 \\ \beta^y & \text{if } y = 1 \text{ \& } s = 0 \\ 1 & \text{if } y = 1 \text{ \& } s = 1. \end{cases} \quad (1)$$

By choosing $\beta^y \neq 0.5$, given a \hat{y} , we can accurately estimate both y and s , which results in a $(\delta^y = 1, \delta^s = 1)$ -HBC classifier. Notice that the classifier in Figure 3 is not a linear classifier anymore, but its capacity is just twice the capacity of the logistic regression. Thus, while keeping the same honesty $\delta^y = 1$, we could improve curiosity from $\delta^s = 0.5$ to $\delta^s = 1$ just by doubling the classifier’s capacity.

The normal mixture is challenging as the range of possible values for \hat{y} is $[0, 1]$. An idea is to define a threshold $\tau' \in [0, 1]$ and divide the range of $[0, 1]$ into four sub-ranges such that:

$$\text{if } \begin{cases} \hat{y} \in [0, \tau') & \text{then we predict } \bar{y} = 0 \text{ \& } \bar{s} = 0 \\ \hat{y} \in [\tau', 0.5) & \text{then we predict } \bar{y} = 0 \text{ \& } \bar{s} = 1 \\ \hat{y} \in [0.5, 1 - \tau') & \text{then we predict } \bar{y} = 1 \text{ \& } \bar{s} = 0 \\ \hat{y} \in [1 - \tau', 1) & \text{then we predict } \bar{y} = 1 \text{ \& } \bar{s} = 1. \end{cases} \quad (2)$$

As we see in Figure 4, a normal mixture cannot guarantee the optimal $(\delta^y = 1, \delta^s = 1)$ -HBC that we could obtain via a hard mixture. Nevertheless, in the following sections, we will show that the idea of dividing the range $[0, 1]$ into four sub-ranges is still useful, especially in white-box situations, where a hard mixture approach is not an option but we can have non-linear classifiers.

Summary. In a black-box view the server can always train two separate classifiers, each with sufficiently high accuracy, and can use the hard mixture of two outputs to build the best achievable

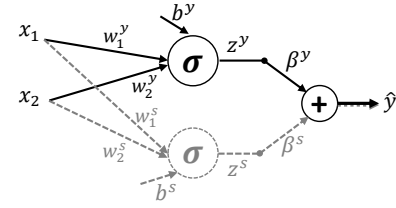


Figure 3: Two logistic regression classifiers, each trained on x separately: z^y for y and z^s for s . Classifiers outputs are mixed with each other such that the final output, \hat{y} , is a real-valued scalar informative about both attributes.

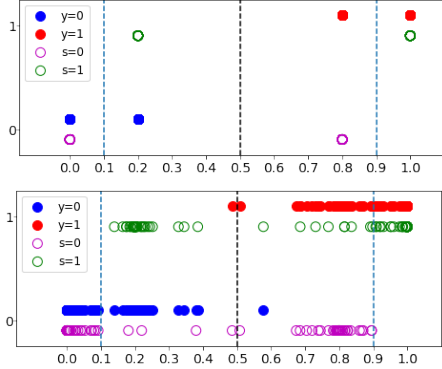


Figure 4: Produced \hat{y} in Figure 3 for the dataset in Figure 2, using $\beta^y = 0.8$ and $\beta^s = 0.2$. (Top) the hard mixture achieves $\delta^y = 1$ and $\delta^s = 1$ (Eq (1)). (Bottom) the normal mixture achieves $\delta^y = .99$ and $\delta^s = .95$ using $\tau' = .1$ (Eq (2)).

(δ^y, δ^s) -HBC classifier; no matter what type of classifier is used. Basically, the server can mostly get the same performance if it could separately run two classifiers on the data. We emphasize that the server’s motivation for such a mixture classifier (and not just simply using two separate classifiers) is that the shape of the classifier’s output \hat{y} depends on Y . Thus, a limitation is that such an attack works only if $S \leq Y$, which in general is not always the case. For instance, let $Y = 5$, $S = 3$, $\beta^y = 0.8$, and $\beta^s = 0.2$. If the hard outputs for y is $\lfloor z^y \rfloor = [0, 0, 0, 1, 0]$ and for s is $z^s = [0, 1, 0]$, then by observing $\hat{y} = [0, 0.2, 0, 0.8, 0]$, server can estimate both y and s while looking very honest. But if $S > Y$, then the server cannot easily encode the private attribute via a mixture model, because we assume that the cardinality of the classifier’s output is limited to Y as it is supposed to look like a standard classifier. Thus, situations with $S > Y$ are more challenging, particularly when users have a white-box view and the server is not free to choose any arbitrary architecture, e.g., it has to train an of-the-shelf classifier. In the following, we focus on the white-box view in a general setting.

4 WHITE-BOX VIEW: A GENERAL SOLUTION

Theoretically, the classifier’s output \hat{y} , as a real-valued vector, can carry an infinite amount of information. Thus, releasing \hat{y} without imposing any particular constraint can reveal any private information and even can be used to (approximately) reconstruct data x . One can imagine a hash function that maps each x to a specific \hat{y} , and consequently, by observing \hat{y} , we can reconstruct x [57]. However, the complexity of real-world data, assumptions on the

required honesty, white-box view, requirement of soft outputs, and other practical constraints will rule out such trivial solutions. Here, we discuss the connection between the curiosity of a classifier and the entropy of its output, and then we formulate the problem of establishing a desired trade-off for an HBC classifier \mathcal{F} and its corresponding attack \mathcal{G} into an information-theoretical framework.

4.1 Curiosity and Entropy

The entropy of a random variable x is the expected value of the information content of that variable; also called self-information: $H(x) = \mathbb{I}(x; x)$. When we are looking for target information in the data, then the presence of other potentially unrelated information in that data could make the extraction of target information more challenging. In principle, such unrelated information would act as noise for our target task. For example, when looking for a target attribute, a trained DNN classifier takes data \mathbf{x} (usually with very high entropy) and produces a probability distribution over the possible outcomes $\hat{y} \in [Y]$ with much lower entropy compared to \mathbf{x} . Although \hat{y} contains much less information than \mathbf{x} in the sense that it has less entropy, it is considered more informative *w.r.t.* the target y .

On the other hand, there is a relationship between the entropy of a classifier’s output, $H(\hat{y}) \in [0, \log(Y)]$, and the curiosity δ^s of an attack \mathcal{G} . The larger $H(\hat{y})$, the more information is carried by \hat{y} , thus the higher the chance to reveal information unrelated to the target task. For example, assume that $Y = 4$, $y = 1$, and y is independent of s . In the extreme case when the classifier’s output is $\hat{y} = [0, 1, 0, 0]$ (that means $H(\hat{y}) = 0$), then \hat{y} carries no information about s and adding any information about s would require increasing the entropy of the output \hat{y} .

In supervised learning, the common loss function for training DNN classifiers is *cross entropy*: $H_y(\hat{y}) = -\sum_{i=0}^{Y-1} \mathbb{I}_{(y=i)} \log \hat{y}_i$; that inherently minimizes $H(\hat{y})$ during training. However, since data is usually noisy, we cannot put any upper bound on $H_y(\hat{y})$ at inference time. In practice, minimizing $H(\hat{y})$, alongside $H_y(\hat{y})$, might help in keeping $H(\hat{y})$ low at inference time, which turns out to be useful for some applications like semi-supervised learning [18]. But there is no guarantee that a classifier will always produce a minimum- or bounded-entropy output at inference time. This fact somehow serves as the main motivation of our work for encoding private attributes of the classifier’s input into the classifier’s output; explained in the following two attacks.

4.2 Regularized Attack

We first introduce a method, for training any classifier to be HBC, in situations where sensitive attribute is binary ($S = 2$) and the server only has access to the soft output ($\hat{y} \in \mathbb{P}^Y$); see Figure 1. The idea is to enforce classifier \mathcal{F} to explicitly encode s into the entropy of \hat{y} by *regularizing* the loss function on \mathcal{F} . In general, there are two properties of \hat{y} that one can utilize for creating an HBC classifier:

1. **Argmax**: as the usual practice, we use the index of the maximum element in \hat{y} to predict y . This helps the classifier to satisfy the honesty requirement.

2. **Entropy**: the entropy of \hat{y} can have at least two states: (i) be close to the maximum entropy, *i.e.*, $H(\hat{y}) = \log Y$, or (ii) be close to the minimum entropy, *i.e.*, $H(\hat{y}) = 0$.

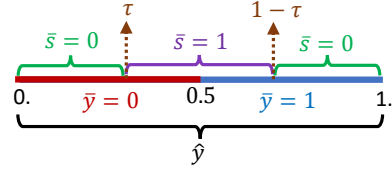


Figure 5: Observing $\hat{y} \in [0, 1]$, we can predict $\bar{y} = 0$ if $\hat{y} < .5$, otherwise $\bar{y} = 1$. We train the classifier such that for samples with $s = 0$, \hat{y} gets close to the borders (0 or 1) depending on y ; otherwise, for $s = 1$, \hat{y} gets far from the borders. Using a threshold τ , we predict binary attributes y and s from \hat{y} .

We show that $H(\hat{y})$ can be used for predicting a binary s , while not interfering with the $\text{argmax}(\hat{y})$ that is preserved for y . Without loss of generality, let us assume $Y = 2$. Consider observing $\hat{y} \in \mathbb{P}$; that is equivalent to $\hat{y} = [\hat{y}_0, \hat{y}_1] \in \mathbb{P}^2$, when $\hat{y}_1 \equiv \hat{y}$ and $\hat{y}_0 = 1 - \hat{y}_1$. Figure 5 shows how we can use the single real-valued \hat{y} to predict two attributes. For example, $\hat{y} = [.95, .05]$ and $\hat{y} = [.75, .25]$ have the same argmax but different entropies: 0.29 and 0.81, respectively.

Training. Choosing any arbitrary classifier \mathcal{F} , the server can train \mathcal{F} with the following loss function:

$$\begin{aligned} \mathcal{L}^b &= \beta^y H_y(\hat{y}) + \beta^s (\mathbb{I}_{(s=0)} - \mathbb{I}_{(s=1)}) H(\hat{y}) = \\ & \beta^y \left(-\sum_{i=0}^{Y-1} \mathbb{I}_{(y=i)} \log \hat{y}_i \right) + \beta^s (\mathbb{I}_{(s=0)} - \mathbb{I}_{(s=1)}) \left(-\sum_{i=0}^{Y-1} \hat{y}_i \log \hat{y}_i \right), \end{aligned} \quad (3)$$

where multipliers β^y and β^s aim to control the trade-off between honesty and curiosity. In Eq (3), in the first term, we have the cross-entropy and in the second term, we have Shannon entropy that aims to minimize the entropy of \hat{y} for samples of $s = 0$, while maximizing the entropy of \hat{y} for samples of $s = 1$.

Attack. At inference time, when the server observes \hat{y} , it computes $H(\hat{y})$, and using a threshold $\tau \in [0, 1]$, estimates \bar{s} :

$$\bar{s} = \mathcal{G}(\hat{y}) = \begin{cases} 0, & \text{if } H(\hat{y}) \leq \tau \\ 1, & \text{otherwise.} \end{cases} \quad (4)$$

Thus, the attack \mathcal{G} is a simple threshold function, and τ is optimized using the validation set during training, as we explain in Section 5.

4.3 Parameterized Attack

In this section, we present our general solution that works for $S \geq 2$, and for both raw and soft outputs; see Figure 1.

4.3.1 An Information Bottleneck Formulation. Remember the Markov chain $(y, s) \rightarrow \mathbf{x} \rightarrow \hat{y} \rightarrow \hat{s}$. We assume that the server is constrained to a specific family of classifiers \mathbb{F} , *e.g.*, a specific DNN architecture that has to be as honest as a standard classifier. The server looks for a $\mathcal{F}^* \in \mathbb{F}$ that maps the users’ data \mathbf{x} into a vector \hat{y} such that \hat{y} is as informative about both y and s as possible.

Formally, \mathcal{F}^* can be defined as the solution of the following mathematical optimization:

$$\min_{(\mathbf{x}, y, s) \leftarrow \mathcal{D}, \mathcal{F} \in \mathbb{F}, \hat{y} \leftarrow \mathcal{F}(\mathbf{x})} [I = \beta^x \mathbb{I}(\hat{y}; \mathbf{x}) - \beta^y \mathbb{I}(\hat{y}; y) - \beta^s \mathbb{I}(\hat{y}; s)], \quad (5)$$

where β^x , β^y , and β^s are Lagrange multipliers that allow us to move along different possible local minimas and all are non-negative real-valued². Eq (5) is an extension of the *information bottleneck* (IB) formulation [71], where the optimal \hat{y} , produced by \mathcal{F}^* , is decided based on its relation to the three variables, x , y , and s . By varying the β multipliers, we can explore the trade-off between *compression* at various rates, *i.e.*, by minimizing $I(\hat{y}; x)$, and the amount of information we aim to preserve, *i.e.*, by maximizing $I(\hat{y}; y)$ and $I(\hat{y}; s)$. Particularly for DNNs, it is shown that compression might help the classifier to achieve better generalization [72].

An Intuition. For better understanding, assume that for $\beta^s = 0$ (*i.e.*, the standard classifiers) the optimal solution \mathcal{F}^* obtains a specific value for $\mathcal{I} = \mathcal{I}^*$ in Eq (5). Now, assume that the server aims to find an HBC solution, by setting $\beta^s > 0$. Then, in order to maintain the same value of \mathcal{I}^* with the same MI for the target attribute $I(\hat{y}; y)$, the term $I(\hat{y}; x)$ must be increased, because $I(\hat{y}; s) \geq 0$ if $\beta^s > 0$. Since for deterministic classifiers $I(\hat{y}; x) = H(\hat{y})$, we conclude that when the server wants to encode information about both y and s in the output \hat{y} , then the output's entropy must be higher than if it was only encoding information about y ; which shows that our formulation in Eq (5) is consistent with our motivation for exploiting the capacity of $H(\hat{y})$. In Section 5, we provide more intuition on this through some experimental results.

4.3.2 Variational Estimation. Now, we analyze a server that is computationally bounded and has only access to a sample of the true population (*i.e.*, a training dataset \mathbb{D}^{train}) and wants to solve Eq (5) to create a (near-optimal) HBC classifier \mathcal{F} . We have

$$\begin{aligned} \mathcal{I} &= \beta^x I(\hat{y}; x) - \beta^y I(\hat{y}; y) - \beta^s I(\hat{y}; s) \\ &= \beta^x H(\hat{y}) - \beta^x H(\hat{y}|x) - \beta^y H(y) + \beta^y H(y|\hat{y}) - \beta^s H(s) + \beta^s H(s|\hat{y}). \end{aligned}$$

Since for a fixed training dataset $H(y)$ and $H(s)$ are constant during the optimization and for a deterministic \mathcal{F} we have $H(\hat{y}|x) = 0$, we can simplify Eq (5) as

$$\min_{(x, y, s) \leftarrow \mathcal{D}, \mathcal{F} \in \mathbb{F}, \hat{y} \leftarrow \mathcal{F}(x)} [\mathcal{H} = \beta^x H(\hat{y}) + \beta^y H(y|\hat{y}) + \beta^s H(s|\hat{y})]. \quad (6)$$

Eq (6) can be interpreted as an optimization problem that aims to minimize the entropy of \hat{y} subject to encoding as much information as possible about y and s into \hat{y} . Thus, the optimization seeks for a function \mathcal{F}^* to produce a low-entropy \hat{y} such that \hat{y} is only informative about y and s and no information about anything else. Multipliers β^y and β^s specify how y and s can compete with each other for the remaining capacity in the entropy of \hat{y} ; that is challenging, particularly, when y and s are independent.

Different constraints on the server can lead to different optimal models. As we observed, in a black-box view with $S \leq Y$ and arbitrary \mathbb{F} , a solution is achieved by training two separate classifiers with a cross-entropy loss function and an entropy minimization regularizer [18]. First, using stochastic gradient descent (SGD), we train classifier \mathcal{F}^y by setting $\beta^y = 1$ and $\beta^s = 0$ in Eq (6). Second, we train classifier \mathcal{F}^s by setting $\beta^y = 0$ and $\beta^s = 1$. Finally, we build $\mathcal{F} = \beta^y [\mathcal{F}^y] + \beta^s [\mathcal{F}^s]$ as the desired HBC classifier for any choice of $\beta^y \in [0, 1]$ and $\beta^s = 1 - \beta^y$. Notice that the desired value for β^x can be chosen through a cross-validation process.

²Mathematically speaking, we only need two Lagrange multipliers as β^y and β^s are dependent. Here we use a redundant multiplier for the ease of presentation.

Thus, let us focus on the white-box view with a constrained \mathbb{F} , *e.g.*, where the server is required to train a known off-the-shelf classifier. Here, we use the cross-entropy loss function for the target attribute y and train the classifier \mathcal{F} using SGD. However, besides this cross-entropy loss, we also need to look for another loss function for the attribute s . Thus, we need a method to simulate such a loss function for s .

Let $p_{s|\hat{y}}$ denote the true, but unknown, probability distribution of s given \hat{y} , and $q_{s|\hat{y}}$ denote an approximation of $p_{s|\hat{y}}$. Considering the cross entropy between these two distributions $H_{p_{s|\hat{y}}}(q_{s|\hat{y}})$, it is known [41] that

$$H(p_{s|\hat{y}}) = - \sum_{i=0}^{S-1} p_{s|\hat{y}} \log(p_{s|\hat{y}}) \leq H_{p_{s|\hat{y}}}(q_{s|\hat{y}}) = - \sum_{i=0}^{S-1} p_{s|\hat{y}} \log(q_{s|\hat{y}}). \quad (7)$$

This inequality tell us that the cross entropy between the unknown distribution, *i.e.*, $p_{s|\hat{y}}$, and any estimation of it, *i.e.*, $q_{s|\hat{y}}$, is an upper-bound on $H(p_{s|\hat{y}})$; and the equality holds when $q_{s|\hat{y}} = p_{s|\hat{y}}$. Thus, if we find a useful model for $q_{s|\hat{y}}$, then the problem of minimizing $H(p_{s|\hat{y}})$ can be solved through minimization of $H_{p_{s|\hat{y}}}(q_{s|\hat{y}})$.

Training. In practice, parameterized models such as neural networks, are suitable candidates for $q_{s|\hat{y}}$ [41, 55]. For N i.i.d. samples of pairs $\{(\hat{y}^1, s^1), \dots, (\hat{y}^N, s^N)\}$ that represent $p_{s|\hat{y}}$ and are generated via our current classifier \mathcal{F} on a dataset $\mathbb{D}^{train} \sim \mathcal{D}$, we can estimate $H_{p_{s|\hat{y}}}(q_{s|\hat{y}})$ using the empirical cross-entropy as

$$\hat{H}_{p_{s|\hat{y}}}^N(q_{s|\hat{y}}) = - \frac{1}{N} \sum_{n=1}^N \sum_{i=0}^{S-1} \mathbb{I}_{(s=i)} \log(q_{s|\hat{y}}(\hat{y}^n)). \quad (8)$$

Therefore, after initializing a *parameterized* model for $q_{s|\hat{y}}$ to estimate $H(p_{s|\hat{y}})$, we run optimization

$$\min_{q_{s|\hat{y}}} \hat{H}_{p_{s|\hat{y}}}^N(q_{s|\hat{y}}), \quad (9)$$

where we iteratively sample N pairs (s^i, \hat{y}^i) , compute Eq (8), and update the model $q_{s|\hat{y}}$; using SGD. The server will use this additional model $q_{s|\hat{y}}$ for s , alongside the cross entropy loss function for y , to solve Eq (6) for finding the optimal \mathcal{F} . Considering $q_{s|\hat{y}}$ as the attack \mathcal{G} and setting our desired β multipliers, the variational approximation of our general optimization problem in Eq (6) is written as

$$\min_{(x, y, s) \leftarrow \mathcal{D}, \mathcal{F} \in \mathbb{F}, \hat{y} \leftarrow \mathcal{F}(x), \hat{s} \in \mathcal{G}(\hat{y})} [\mathcal{H} = \beta^x \hat{H}(\hat{y}) + \beta^y \hat{H}(y|\hat{y}) + \beta^s \hat{H}(s|\hat{y})], \quad (10)$$

where the joint minimization is performed over both parameterized models \mathcal{F} and \mathcal{G} . Here, $\hat{H}(\cdot)$ and $\hat{H}(\cdot|\cdot)$ denote the empirical entropy and conditional-entropy computed on every sampled batch of data, respectively. A schematic view of this approach is depicted in Figure 6. Using a training dataset $(x, y, s) \in \mathbb{D}^{train} \sim \mathcal{D}$ and in an iterative process, both classifier \mathcal{F} , with loss function $\mathcal{L}^{\mathcal{F}} = \beta^x \hat{H}(\hat{y}) + \beta^y \hat{H}(y|\hat{y}) + \beta^s \hat{H}(s|\hat{y})$, and attack \mathcal{G} , with loss function $\mathcal{L}^{\mathcal{G}} = \hat{H}(s|\hat{y})$, are simultaneously trained. Algorithm 1, in Appendix B, shows the details of our proposed training method.

Summary. While our *regularized* attack (Section 4.2) works by only modifying the loss function of \mathcal{F} , in our *parameterized* attack (this section) we not only need to modify the loss function of \mathcal{F} but also to utilize an additional model \mathcal{G} (*e.g.*, a multi-layer perception)

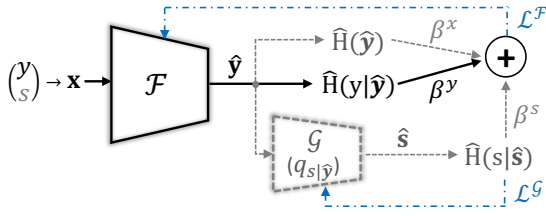


Figure 6: The server iteratively trains the chosen classifier \mathcal{F} and attack \mathcal{G} that is a *parameterized* model. While \mathcal{G} is optimized only based on $\hat{H}(s|\hat{s})$, classifier \mathcal{F} is optimized based on a weighted combination of $\hat{H}(\hat{y})$, $\hat{H}(y|\hat{y})$ and $\hat{H}(s|\hat{s})$ in Eq (10).

to estimate the sensitive attribute. Hence, \mathcal{F} can be seen as a model that aims to perform two roles at the user’s side: (1) as an honest classifier, to estimate the target attribute, and (2) as a curious encoder to encode the sensitive attribute. At the server side, \mathcal{G} acts as a decoder for decoding the received output for estimating the sensitive attribute. Notice that, since the sensitive attribute is a categorical variable, we cannot directly reconstruct s , e.g., using mean squared error in traditional autoencoders, and we need Equation (9).

5 EVALUATION

We present results in several tasks from two real-world datasets and discuss the honesty-curiosity trade-off for different attributes.

5.1 Experimental Setup

5.1.1 Settings. We define a *task* as training a specific classifier \mathcal{F} on a training dataset including samples with two attributes, $y \in [Y]$ and $s \in [S]$, and evaluating \mathcal{F} on a test dataset by measuring honesty δ^y and curiosity δ^s via the attack \mathcal{G} (see Section 2). We have two types of attacks: (i) *regularized*, as explained in Section 4.2, and (ii) *parameterized*, as explained in Section 4.3. For each task there are 3 different scenarios: (1) Std: where \mathcal{F} is a standard classifier without intended curiosity, (2) RawHBC: where \mathcal{F} is trained to be HBC and \mathcal{G} has access to raw outputs $\hat{y} \in \mathbb{R}^Y$, and (3) SoftHBC: where \mathcal{F} is HBC but \mathcal{G} has access only to soft outputs $\hat{y} \in \mathbb{P}^Y$ (see Figure 1). We run each experiment five times, and report mean and standard deviation. For each experiment in SoftHBC, $\tau \in [0, 1]$ in Eq (4) is chosen based on the validation set and is used to evaluate the result on the test set.

5.1.2 CelebA Dataset [35]. This is a dataset including more than 200K celebrity face images, each with 40 binary attributes, e.g., the ‘Smiling’ attribute with values 0:Non-Smile or 1:Smile. We choose attributes that are almost balanced, meaning that there are at least 30% and at most 70% samples for that attribute with value 1. Our chosen attributes are: Attractive, BlackHair, BlondHair, BrownHair, HeavyMakeup, Male, MouthOpen, Smiling, and WavyHair. CelebA is already split into separate training, validation, and test sets. We use the resampled images of size 64×64 . We elaborate more on CelebA and show some samples of this dataset in Appendix F.1.

5.1.3 UTKFace Dataset [86]. This is a dataset including 23,705 face images annotated with attributes of Gender (Male or Female), Race (White, Black, Asian, Indian, or others), and Age (0-116). We use the resampled images of size 64×64 , and randomly split UTKFace into

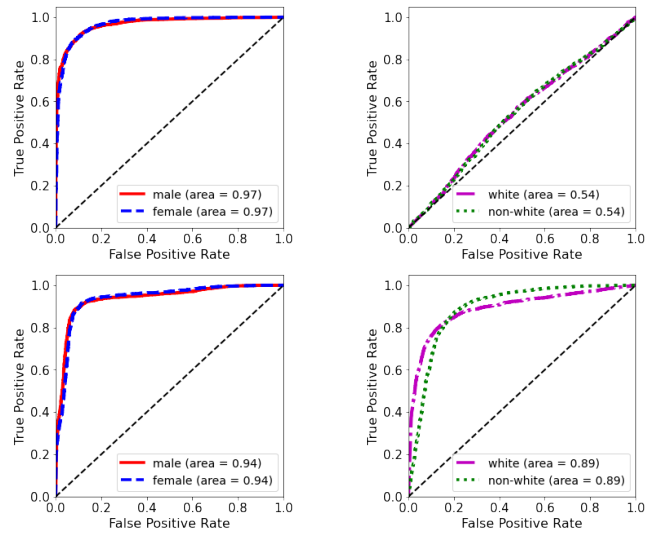


Figure 7: The ROC curve of (top) a standard classifier and (bottom) an HBC classifier trained by the regularized attack in SoftHBC. The target attribute is Gender and the sensitive attribute is Race (in UTKFace), and (β^y, β^s) in Eq (3) is $(.5, .5)$.

subsets of sizes 18964 (80%) and 4741 (20%) for training and test sets, respectively. A subset of 1896 images (10%) from the training set is randomly chosen as the validation set for training. See Appendix F.2 for the details of tasks we define on UTKFace, and some samples.

5.1.4 Architectures. For \mathcal{F} , we use a DNN architecture similar to the original paper of UTKFace dataset [86] that includes 4 convolutional layers and 2 fully-connected layers with about 250K trainable parameters. For \mathcal{G} , we use a simple 3-layer fully-connected classifier with about 2K to 4K trainable parameters; depending on the value of Y . The implementation details for \mathcal{F} and \mathcal{G} are presented in Appendix G. For all experiments, we use a batch size of 100 images, and Adam optimizer [29] with learning rate .001. After fixing β multipliers, we run training for 50 epochs, and choose models of the epoch that both \mathcal{F} and \mathcal{G} achieve the best trade-off for both y and s (based on β^y and β^s) on the validation set, respectively. That is, the models that give us the largest $\beta^y \delta^y + \beta^s \delta^s$ on the validation set during training. Notice that, the fine-tuning is a task at training time, and a server with enough data and computational power can find near-optimal values for (β^y, β^s) , as we do here using the validation set. In the following, all reported values for honesty δ^y and curiosity δ^s of HBC classifiers are the accuracy of the final \mathcal{F} and \mathcal{G} on the test set. Finally, in all Std settings, values in *italic* show the effect of overlearning [66], that is the accuracy of a parameterized \mathcal{G} in inferring a sensitive attribute from a standard classifier.

5.2 UTKFace: Gender vs. Race

As the first result, we present a simple result on UTKFace. We set Gender as the target attribute y and Race (White, Non-White) as the sensitive attribute s . The training set includes 52% Male and 48% Female, where 42% of samples are labeled as White and 58% as Non-White. Figure 7 shows the ROC curves for this experiment. For honesty, in the top-left plot, the standard classifier achieves

Table 1: (A) The characteristics of four attributes vs. Smiling attribute in the training set of CelebA dataset. The honesty (δ^y) and curiosity (δ^s) of (B) standard classifiers when releasing either raw or soft outputs, as well as our proposed HBC classifiers trained for (C) regularized attacks and (D) parameterized attacks. Here, $Y = S = 2$ and $\beta^x = 0$. Values are in percentage (%).

(A) For each attribute: the empirical joint probability and mutual information (MI) with Smiling, besides accuracy of the standard classifier (δ^y)																
MouthOpen				Male				HeavyMakeup				WavyHair				
	Non-Smile	Smile	sum		Non-Smile	Smile	sum		Non-Smile	Smile	sum		Non-Smile	Smile	sum	
0	.386	.122	.508	0	.292	.375	.667	0	.305	.184	.489	0	.322	.166	.488	
1	.103	.389	.492	1	.197	.136	.333	1	.226	.285	.511	1	.303	.209	.512	
sum	.489	.510		sum	.489	.511		sum	.531	.469		sum	.625	.375		
	MI: .231	δ^y : 93.42 \pm .07		MI: .015	δ^y : 97.67 \pm .04			MI: .024	δ^y : 88.89 \pm .19			MI: .003	δ^y : 77.29 \pm .55			
Setting	(β^y, β^s)		s: MouthOpen		s: Male		s: HeavyMakeup		s: WavyHair							
			δ^y	δ^s	δ^y	δ^s	δ^y	δ^s	δ^y	δ^s	δ^y	δ^s				
(B) Overlearning [66]																
Std	raw	(1., 0.)	92.15 \pm .04	79.56 \pm .32	92.15 \pm .04	68.20 \pm .03	92.15 \pm .04	60.96 \pm .10	92.15 \pm .04	57.93 \pm .21						
	soft			79.22 \pm .14		68.20 \pm .00		60.53 \pm .11		57.90 \pm .00						
(C) Regularized Attack																
SoftHBC		(.7, .3)	91.90 \pm .01	84.73 \pm .50	91.78 \pm .15	90.14 \pm .37	92.03 \pm .12	80.79 \pm .47	92.11 \pm .15	68.06 \pm .72						
		(.5, .5)	91.83 \pm .13	89.08 \pm .04	91.65 \pm .10	94.20 \pm .18	91.87 \pm .03	85.27 \pm .15	91.98 \pm .18	72.36 \pm .83						
		(.3, .7)	91.75 \pm .10	91.22 \pm .25	91.60 \pm .11	96.02 \pm .09	91.58 \pm .10	87.09 \pm .43	91.73 \pm .07	73.52 \pm .77						
(D) Parameterized Attack																
RawHBC		(.7, .3)	91.84 \pm .07	93.40 \pm .09	92.36 \pm .02	97.21 \pm .10	92.12 \pm .09	88.63 \pm .04	92.17 \pm .04	76.74 \pm .37						
		(.7, .3)	91.49 \pm .44	85.94 \pm .76	91.72 \pm .17	79.39 \pm 2.4	91.84 \pm .10	84.45 \pm .15	92.12 \pm .09	57.90 \pm .00						
SoftHBC		(.5, .5)	90.20 \pm 1.1	88.91 \pm .83	90.17 \pm .66	92.73 \pm .76	91.49 \pm .11	86.19 \pm .18	92.20 \pm .12	61.51 \pm .77						
		(.3, .7)	82.55 \pm .27	93.15 \pm .07	67.65 \pm 6.4	97.41 \pm .06	70.75 \pm 4.3	88.62 \pm .52	89.27 \pm 1.9	67.57 \pm 3.7						

0.97 area under the ROC curve (AUC), whereas in the bottom-left plot the classifier is HBC but it still achieves a considerable 0.94 AUC. For curiosity, the standard classifier in the top-right plot is not informative about Race and it basically is as good as a random guess. But, the HBC classifier in the bottom-right can achieve 0.89 AUC on predicting Race via the regularized attack.

5.3 CelebA: Smiling vs. Other Attributes

Here, we consider a smile detection task on CelebA where y is set to Smiling and s is set to one of MouthOpen, Male, HeavyMakeup, or WavyHair. To better understand the results, Table 1 (A) shows the characteristics of these attributes compared to Smiling attribute in the training dataset. As expected, MouthOpen is the most correlated one (empirical MI is .231 bits) while WavyHair is the least correlated attribute (empirical MI is .003 bits), to Smiling. We also show the test accuracy of \mathcal{F} as a standard classifier trained for each attribute. We see that Male is the easiest attribute (97% accuracy) and WavyHair is the most difficult one (77% accuracy). Table 1 (B), (C), and (D) show how these factors, *correlation* and *easiness*, affect the performance of HBC classifiers: for different attacks and trade-offs based on β^y and β^s . Our findings are:

1. The outputs of a standard binary classifier, either raw or soft, do not reveal information about sensitive attributes more than what one could already infer from the knowledge about underlying data distribution. Thus, while overlearning [66] has shown serious problems when the server observes a subset of internal representations, here we see that when the server only observes the output of a binary classifier, then overlearning is not a major problem (we show this in other settings as well). On the other hand, HBC classifiers can

effectively learn to encode the sensitive attribute in a single-valued output (that is more restricted than the internal representation).

2. In regularized attacks in Table 1 (B), with a very small loss in δ^y (less than 1%), we can get δ^s very close to the accuracy that we could have achieved if we could run a separate classifier for the sensitive attribute. In parameterized attacks, it is easier to encode sensitive information into the raw output than the soft output. The attacks in RawHBC are highly successful in curiosity in all four cases in Table 1 (C), almost without any damage to the honesty. On the other hand, in SoftHBC it is harder to establish an efficient trade-off between honesty and curiosity. Since $\text{softmax}(x) = \text{softmax}(x + a)$ for all $a \in \mathbb{R}$, there are infinitely many vectors in \mathbb{R}^Y (in RawHBC) that can be mapped into the same vector in \mathbb{P}^Y (in SoftHBC). However, what one can learn about the sensitive attribute in SoftHBC is still much more than Std. In the following, we will see that when $Y > 2$, parameterized attack in SoftHBC is also very successful.

3. The easiness of the sensitive attribute plays an important role. For example, in face image processing, gender classification is in general an easier task than detecting heavy makeup (as it might be easier for human beings as well). Therefore, while MI between Smiling and HeavyMakeup is larger than Smiling and Male, the curiosity in inferring Male attribute is more successful than HeavyMakeup. Moreover, while (due to MI) δ^s of MouthOpen is about 11% more than Male in Std, in contrast to overlearning attack, the correlation is not that important in HBC settings compared to the easiness, as we see that for Male all attacks are as successful as MouthOpen. For the same reason, for attribute WavyHair it is more difficult to achieve high curiosity as it is not an easy task. It is worth noting that, for difficult attributes we may be able to improve the curiosity by optimizing \mathcal{F} for that specific attribute, e.g.,

Table 2: Results where the target attribute is HairColor and the sensitive attribute is (A) Male, (B) Smiling, or (C) Attractive. Here $Y=3$, $S=2$, and $\beta^x=0$. Values are in percentage (%).

Setting	Attack	(β^y, β^s)	δ^y	δ^s
(A) Male (class distribution 0:66%, 1:34%)				
Std	raw [66]	(1., .0)	92.94 ± .58	75.86 ± .03
	Parameterized	(.7, .3)	92.85 ± .42	97.27 ± .09
RawHBC	Parameterized	(.7, .3)	92.12 ± .27	94.23 ± 1.4
		(.5, .5)	91.98 ± .12	96.87 ± .13
	SoftHBC	(.3, .7)	67.63 ± .25	97.52 ± .08
		(.7, .3)	92.79 ± .22	95.11 ± .22
	Regularized	(.5, .5)	92.78 ± .28	96.68 ± .12
		(.3, .7)	92.96 ± .33	97.02 ± .04
(B) Smiling (class distribution 0:49%, 1:51%)				
Std	raw [66]	(1., .0)	92.94 ± .58	56.12 ± .21
	Parameterized	(.7, .3)	92.56 ± .46	91.79 ± .12
RawHBC	Parameterized	(.7, .3)	92.96 ± .44	88.59 ± 1.8
		(.5, .5)	91.27 ± .35	89.18 ± .06
	SoftHBC	(.3, .7)	88.05 ± 1.3	91.07 ± .51
		(.7, .3)	92.80 ± .28	86.22 ± .38
	Regularized	(.5, .5)	92.52 ± .43	90.67 ± .15
		(.3, .7)	92.31 ± .32	91.59 ± .07
(C) Attractive (class distribution 0:40%, 1:60%)				
Std	raw [66]	(1., .0)	92.94 ± .58	60.73 ± .21
	Parameterized	(.7, .3)	92.81 ± .16	76.97 ± .07
RawHBC	Parameterized	(.7, .3)	93.03 ± .33	71.54 ± .60
		(.5, .5)	92.62 ± .64	76.18 ± .12
	SoftHBC	(.3, .7)	92.66 ± .33	75.57 ± .37
		(.7, .3)	91.95 ± .35	69.06 ± .39
	Regularized	(.5, .5)	92.33 ± .44	74.38 ± .13
		(.3, .7)	91.09 ± .46	75.70 ± .27

through neural architecture search. But for fair comparisons, we use the same DNN for all experiments and we leave the architectures optimization for HBC classifiers to future studies.

5.4 CelebA: HairColor vs. Other Attributes

Moving beyond binary classifiers, in Table 2 we present a use-case of a three-class classifier ($Y=3$) for the target attribute of HairColor, where there are 40% samples of BlackHair, 34% BrownHair, and 26% BlondHair. We consider three sensitive attributes with different degrees of easiness: (A) Male, (B) Smiling, and (C) Attractive. While in Std setting, attacks on overlearning [66] are not very successful (even when releasing the raw outputs), our parameterized attacks, in RawHBC, are very successful without any meaningful damage to the honesty of classifier. In SoftHBC, while it is again harder to train a parameterized attack as successful as RawHBC, we do find successful trade-offs if we fine-tune (β^y, β^s) ; particularly for the regularized attack.

5.5 UTKFace: Sensitive Attributes with $S>2$

To evaluate tasks with $S > 2$, we provide results of several experiments performed on UTKFace in Table 3 and Table 4 (also some

complementary results in Appendix C, Table 7, Tables 8, and Table 9). In each experiment, we set one of Gender, Age, or Race, as y and another one as s , and compare the achieved δ^y and δ^s . See Appendix F.2 for the details of how we created labels for different values of Y and S . Our findings are:

1. An HBC classifier can be as honest as a standard classifier while also achieving a considerable curiosity. For all RawHBC cases, the δ^y of an HBC classifier is very close to δ^y of a corresponding classifier in Std. Moreover, we see that in some situations, making a classifier HBC even helps in achieving a better generalization and consequently getting a slightly better honesty; which is very important as an HBC classifier can look as honest as possible (we elaborate more on the cause of this observation in Appendix D).

2. When having access to raw outputs, the attack is highly successful in all tasks, and in many cases, we can achieve similar accuracy to a situation where we could train \mathcal{F} for that specific sensitive attribute. For example, in Table 3 for $S = 3$ and $Y > 2$, we can achieve about 83% curiosity in inferring the Race attribute from a classifier trained for Age attribute. When looking at Table 8 where Race is the target attribute, we see that the best accuracy a standard classifier can achieve for Race classification is about 85%.

3. In SoftHBC, it is more challenging to achieve a high curiosity via a parameterized attack, unless we sacrifice more honesty. Particularly for tasks with $S > Y$, where the sensitive attribute is more granular than the target attribute. Also, while we observed successful regularized attacks for SoftHBC in Table 1 (C), regularized attacks cannot be applied to tasks with $S > 2$. Yet, even in this case of having only access to soft outputs, the curiosity is much higher than what can be learned from the raw output of a standard classifier (via overlearning attack). Although the curiosity in RawHBC is more successful than SoftHBC, the difference between these two gets smaller as the size of output Y gets larger.

4. Attacks are highly successful when $S \leq Y$, as there is more capacity in the released output. However, the attack is successful in scenarios when $S > Y$ as well. The most challenging case is where $Y = 2$ and when we only have access to the soft outputs, because in these tasks we only release one value (*i.e.*, $\hat{y}_1 = 1 - \hat{y}_2$). Moreover, we see in Table 4 that for SoftHBC with $S = Y = 2$, regularized attacks achieve much better trade-offs than parameterized attacks.

5.6 Entropy Minimization with $\beta^x > 0$

We examine the entropy minimization (*i.e.*, compression) of the classifier’s outputs and its effect on the achieved trade-off for honesty and curiosity. Table 5 compares the results of different values chosen for β^x in Eq (10) (see Figure 6). An important observation is that compression is mostly helpful to the honesty. This is expected, due to our discussion in Section 4.3, and findings in previous related works [72, 75]. Moreover, we see that compression is more effective in improving the honesty of the classifier in situations where we assign more weight to the curiosity of the classifier.

Although Table 5 shows that large compression hurts curiosity more, this is another trade-off that a server can utilize to make the HBC classifier less suspicious. It is important to observe that the average entropy of the classifier’s output, shown by $\mathbb{H}(\hat{y})$, is directly related to the curiosity weight β^s . The more curious a classifier is, the larger will be the entropy of the output. In Figure 8, we plot

Table 3: The honesty (δ^y) and curiosity (δ^s) of an HBC classifier trained via parameterized attack, where the target attribute (y) is Age and the sensitive attribute (s) is Race. Values are in percentage (%)

		$S = 2$			$S = 3$			$S = 4$			$S = 5$		
		(β^y, β^s)	δ^y	δ^s	δ^y	δ^s	δ^y	δ^s	δ^y	δ^s	δ^y	δ^s	
$Y = 2$	NC	(1.,0.)	83.03 ± .33	62.44 ± .97	83.03 ± .33	52.34 ± 1.4	83.03 ± .33	45.74 ± .13	83.03 ± .33	44.67 ± .17			
	RawHBC	(.7,.3)	83.60 ± .22	86.77 ± .11	83.22 ± .27	80.06 ± .88	83.10 ± .18	74.12 ± .33	83.07 ± .13	72.52 ± .34			
	SoftHBC	(.7,.3)	82.87 ± .21	69.89 ± .40	82.88 ± .27	63.28 ± .30	81.67 ± .54	55.01 ± .57	81.94 ± .81	53.84 ± .93			
$Y = 3$	NC	(1.,0.)	81.09 ± .37	65.99 ± .40	81.09 ± .37	56.37 ± .47	81.09 ± .37	48.28 ± .49	81.09 ± .37	46.67 ± .41			
	RawHBC	(.7,.3)	81.67 ± .33	86.35 ± .65	81.35 ± .28	83.42 ± .29	81.23 ± .29	76.62 ± .56	81.30 ± .28	76.24 ± .59			
	SoftHBC	(.7,.3)	81.19 ± .02	79.07 ± .15	80.76 ± .36	72.90 ± .23	80.10 ± .32	66.40 ± 1.1	80.29 ± .07	67.74 ± .30			
$Y = 4$	NC	(1.,0.)	68.59 ± .24	66.93 ± .39	68.59 ± .24	58.02 ± .55	68.59 ± .24	49.23 ± .53	68.59 ± .24	41.13 ± .15			
	RawHBC	(.7,.3)	68.59 ± .27	86.91 ± .27	68.59 ± .49	84.17 ± .26	68.29 ± .26	79.46 ± 1.6	68.40 ± .50	78.79 ± .13			
	SoftHBC	(.7,.3)	68.15 ± .42	78.10 ± .32	67.45 ± .25	74.70 ± .43	66.30 ± .22	69.01 ± 1.3	65.70 ± .51	69.37 ± .28			
$Y = 5$	NC	(1.,0.)	62.24 ± .61	67.00 ± .32	62.24 ± .61	58.69 ± .44	67.00 ± .32	50.37 ± .35	67.00 ± .32	49.22 ± .62			
	RawHBC	(.7,.3)	62.72 ± .20	86.63 ± .05	62.46 ± .49	83.79 ± .16	61.84 ± .77	78.53 ± .80	62.40 ± .31	78.28 ± .20			
	SoftHBC	(.7,.3)	62.08 ± .14	79.41 ± .75	61.35 ± .58	74.57 ± .91	60.90 ± .39	69.32 ± .49	60.70 ± .52	69.56 ± 1.2			
		(.5,.5)	58.53 ± .85	86.88 ± .34	56.63 ± .79	84.06 ± .09	53.27 ± .14	77.45 ± .11	54.61 ± .20	77.30 ± 1.0			

Table 4: The honesty (δ^y) and curiosity (δ^s) of an HBC classifier trained via different attacks, where the target attribute (y) is Race and the sensitive attribute (s) is Gender. Values are in percentage (%).

				$Y = 2$		$Y = 3$		$Y = 4$		$Y = 5$	
Setting	Attack	(β^y, β^s)	δ^y	δ^s	δ^y	δ^s	δ^y	δ^s	δ^y	δ^s	
$S = 2$	Std	raw [66]	(1.,0.)	87.67 ± .49	56.11 ± .44	85.50 ± .13	56.04 ± 2.4	81.15 ± .23	58.91 ± .74	80.69 ± .23	61.06 ± 2.0
	RawHBC	Parameterized	(.7,.3)	88.10 ± .22	89.36 ± .03	85.77 ± .30	89.14 ± .21	82.05 ± .14	89.19 ± .33	81.43 ± .34	89.00 ± .42
			(.5,.5)	83.67 ± .27	80.17 ± 2.1	82.89 ± .59	85.69 ± .54	70.36 ± 1.5	86.05 ± 1.6	77.40 ± .35	88.36 ± .45
	SoftHBC	Regularized	(.7,.3)	87.30 ± .29	75.29 ± .89	84.89 ± .20	77.53 ± .33	80.64 ± .31	76.29 ± .45	80.24 ± .09	78.01 ± .62
			(.5,.5)	86.61 ± .15	83.23 ± .27	83.89 ± .44	86.61 ± .61	79.03 ± .21	86.29 ± .32	78.61 ± .65	85.97 ± .28

Table 5: The effect of entropy minimization on the UTKFace dataset for the target attribute Age and the sensitive attribute Race where $Y = S = 3$. We also show the average of the entropy of classifier’s output by $\tilde{H}(\hat{y})$ in bits.

		$\beta^x = .0$			$\beta^x = .2$			$\beta^x = .4$			$\beta^x = .8$		
		δ^y	δ^s	$\tilde{H}(\hat{y})$	δ^y	δ^s	$\tilde{H}(\hat{y})$	δ^y	δ^s	$\tilde{H}(\hat{y})$	δ^y	δ^s	$\tilde{H}(\hat{y})$
Std	(1., .0)	81.43±.37	58.14±.06	.48±.05	81.04±.09	56.13±.14	.45±.06	80.90±.27	58.60±.07	.40±.05	81.15±.41	55.94±.26	.32±.02
	(.7, .3)	81.32±.31	82.97±.27	.63±.05	81.59±.42	82.90±.24	.45±.02	81.62±.37	81.98±.77	.36±.02	81.61±.31	81.40±.41	.28±.03
RawHBC	(.5, .5)	80.23±.42	84.82±.35	.76±.03	80.60±.14	84.60±.32	.54±.03	80.95±.41	83.86±.54	.39±.01	81.26±.22	83.92±.32	.26±.01
	(.7, .3)	81.04±.26	73.34±.29	.72±.02	81.26±.28	69.77±.50	.51±.02	81.01±.55	64.54±2.5	.37±.04	81.01±.46	56.28±.51	.27±.05
SoftHBC	(.5, .5)	69.12±.91	80.81±.97	1.05±.01	76.20±.29	76.27±.37	.76±.01	80.13±.15	74.01±.46	.58±.01	80.59±.16	70.28±.15	.33±.02

the histogram of normalized observations of the output’s entropy for the setting of $Y=S=3$ for five scenarios in Table 5. We see that the output of an HBC classifier tends to have larger entropy than a standard classifier. Moreover, large entropies, *i.e.*, more than 1 bit, are more common in SoftHBC than RawHBC. A reason for this is that the capacity of soft outputs is smaller than raw outputs; hence, the classifier tends to take more advantage of the existing capacity. Interestingly, when we use the entropy minimization with $\beta^x > 0$, then we observe that the entropy distribution for an HBC classifier has a smaller tail (compared to $\beta^x = 0$) and thus can become even less suspicious than the standard classifier when we release the raw output. Finally, in releasing the soft outputs, it is a bit more challenging to keep the average entropy low.

5.7 Pruning HBC Classifiers

As DNN classifiers are mostly overparameterized, a hypothesis might be that HBC classifiers utilize the extra capacity of DNNs for extracting patterns that correspond to the sensitive attribute. Thus, one could say that reducing a DNN’s capacity, using a pruning technique [21], might make it more difficult for a classifier to be curious. For example, a user who does not trust a server can take the classifier \mathcal{F} and perform some pruning technique, before making inferences, hoping that pruning will not damage the honesty but will reduce curiosity. Figure 9 shows the honesty-curiosity trade-off for different amounts of pruning. As the pruning technique, we use

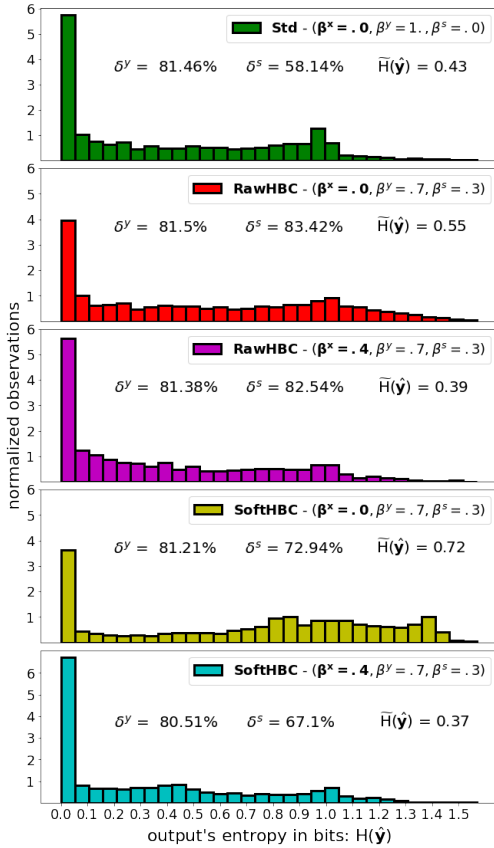


Figure 8: The normalized density histogram (i.e., number of observations) of the output entropies of classifier with and without compression. Here, y is Age and s is Race with $Y = S = 3$ on UTKFace. From top to bottom, classifier is trained (1) standard, (2, 3) RawHBC where we release the raw outputs, and (4, 5) SoftHBC where we release the soft outputs. For (2) and (4) we set $\beta^x = .0$, and for (3) and (5) we set $\beta^x = .4$.

$L1$ -Unstructured implemented in [56], where parameters with the lowest $L1$ -norm are set to zero at inference time (i.e., no re-training).

1. Comparing the top plots (where $\beta^x = 0$) with bottom plots (where $\beta^x = 0.4$) in Figure 9, we see that classifiers without compression show more tolerance to a large amount of pruning (>60%) than classifiers with compression. This might be due to the additional constraint that we put on classifiers with compression.

2. By pruning less than 50% of the parameters, there is no significant drop in the honesty in both standard and HBC classifiers. However, for the curiosity, HBC classifiers show different behaviors in different settings. When there is no compression ($\beta^x = 0$), the curiosity of RawHBC shows faster and much larger drops, compared to SoftHBC. When there is compression ($\beta^x = 0.4$), then drops are not different across these settings.

3. If we prune more than 50% of the parameters, the drop in accuracy for both the target and sensitive attributes are significant. However, SoftHBC settings interestingly show better performances than RawHBC, for both target and sensitive attributes. While RawHBC

Table 6: An HBC teacher can train an HBC student via KD. The target attribute is Smiling in CelebA and the setting is SoftHBC with $(\beta^x, \beta^y, \beta^s) = (.0, .5, .5)$.

	s	classifier	Parameterized		Regularized	
			δ^y	δ^s	δ^y	δ^s
Mouth	Teacher		90.04 ± .61	89.14 ± .68	91.61 ± .24	89.01 ± .34
	Student		90.30 ± .47	88.07 ± .65	91.29 ± .11	83.69 ± .38
Open	Teacher		90.17 ± .66	92.73 ± .76	91.57 ± .16	94.11 ± .21
	Student		90.20 ± .51	91.62 ± .73	91.01 ± .19	85.71 ± .46
Male	Teacher		90.73 ± 1.1	85.17 ± .97	91.56 ± .21	85.58 ± .27
	Student		91.00 ± .43	82.62 ± .49	91.23 ± .15	81.05 ± .15
Heavy	Teacher		91.99 ± .21	61.05 ± 1.8	91.86 ± .08	71.12 ± 1.2
	Student		91.78 ± .17	59.72 ± 1.1	91.60 ± .15	68.61 ± .96
Wavy	Teacher					
	Student					
Hair	Teacher					
	Student					

has usually shown better performance in previous sections, in this specific case SoftHBC performs interestingly very well.

In sum, although Figure 9 shows that pruning can damage curiosity more than honesty in settings without compression. We cannot guarantee that pruning at inference time is an effective defense against HBC models, as adding compression constraint will help the server to make HBC models more tolerable to pruning, and with the amount of pruning up to 50% the curiosity remains high.

5.8 Transferring Curiosity via HBC Teachers

Transferring knowledge from a *teacher* classifier, trained on a large labeled data, to a *student* classifier, that has access only to a (small) unlabeled data, is known as “knowledge distillation” (KD) [20, 22]. Since KD allows to keep sensitive data private, it has found some applications in privacy-preserving ML [52, 77]. Let us consider a user that owns a private unlabeled dataset and wants to train a classifier on this dataset, and a server that provides a teacher classifier trained on a large labeled dataset. A common technique in KD is to force the student to mimic the teacher’s behavior by minimizing the KL-divergence between the teacher’s soft outputs, $\hat{y}^{Teacher}$, and student’s soft outputs, $\hat{y}^{Student}$ [22]:

$$\mathcal{L}^{KL} = \sum_{i=1}^Y \hat{y}_i^{Teacher} \log \left(\frac{\hat{y}_i^{Teacher}}{\hat{y}_i^{Student}} \right). \quad (11)$$

We show that if the teacher is HBC, then the student trained using \mathcal{L}^{KL} can also become HBC. We run experiments, similar to Section 5.3, by assigning 80% of the training set (with both labels) to the server, and 20% of the training set (without any labels) to the user. At the server side, we train an HBC teacher via both regularized and parameterized attacks, and then at the user side, we use the trained teacher classifier to train a student only via \mathcal{L}^{KL} and the user’s unlabeled dataset. We also set the student’s DNN to be half the size of the teacher’s in terms of the number of trainable parameters in each layer. We set $(\beta^y, \beta^s) = (.5, .5)$ and consider the SoftHBC setting, because \mathcal{L}^{KL} is based on mimicking the teacher’s soft outputs.

Table 6 shows that KD works well for the honesty in both cases; showing that an HBC teacher can look very honest in transferring knowledge of the target attribute. For the curiosity, the parameterized attacks transfer the knowledge very well and are usually better

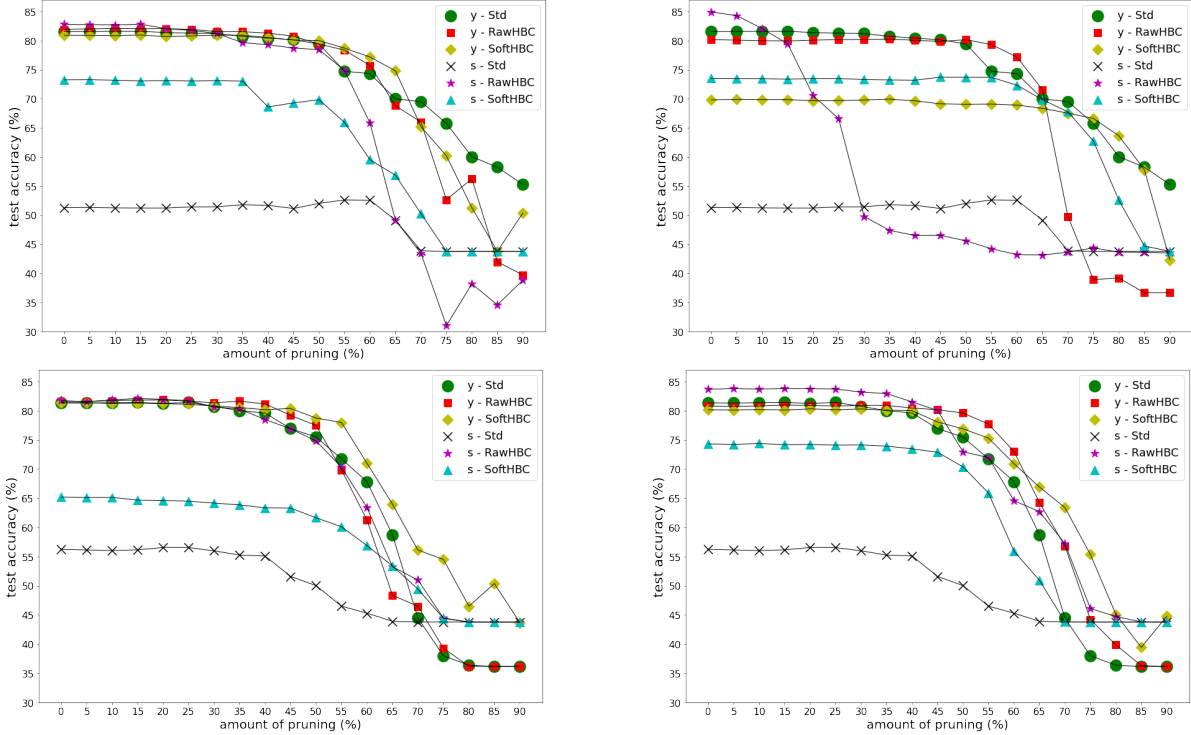


Figure 9: The effect of pruning, with L1 norm, where the classifier \mathcal{F} when y is Age, s Race, and $Y = S = 3$. (Top-Left) $\beta^x, \beta^y, \beta^s = (0, .7, .3)$. (Top-Right) we set $\beta^x, \beta^y, \beta^s = (0, .5, .5)$. (Bottom-Left) $\beta^x, \beta^y, \beta^s = (.4, .7, .3)$. (Bottom-Right) we set $\beta^x, \beta^y, \beta^s = (.4, .5, .5)$.

than regularized attacks. A reason might be that for regularized attacks the chosen threshold τ for the teacher classifier is not the best choice for replicating the attack in its students. Though, for WavyHair as an uncorrelated and difficult attribute, the regularized attack achieves a better result.

Overall, this capability of transferring curiosity shows another risk of HBC classifiers provided by semi-trusted servers, especially that such a teacher-student approach has shown successful applications in semi-supervised learning [24, 58, 68, 84]. We note that there are several aspects, such as the teacher-student data ratio and distribution shift, or different choices of the loss function, that may either improve or mitigate this attack. For example, the entropy of teacher’s output in KD can be controlled via a “temperature” parameter t inside the softmax , i.e., $\hat{y}_i^p = e^{\hat{y}_i^r/t} / \sum_{j=0}^{Y-1} e^{\hat{y}_j^r/t}$, which might further affect the achieved honesty-curiosity. We leave further investigation into these aspects of KD for future studies.

6 DISCUSSIONS AND RELATED WORK

We discuss related work and potential proactive defenses against the vulnerabilities highlighted in this paper.

6.1 Proactive Investigation

The right to information privacy is long known as “the right to select what personal information about me is known to what people” [81]. The surge of applying ML to (almost) all tasks in our everyday lives has brought attention to the ethical aspects of ML [27]. To reconstruct the users’ face images from their recordings of speech [50, 80]

raised the concern of “tying the identity to biology” and categorizing people into gender or sexual orientation groups that they do not fit well [23]. The estimation of ethnicity [76] or detecting sexual orientation [79] from facial images has risen concerns about misusing such ML models by adversaries that seek to determine people of minority groups. It is shown in [49] that a widely used healthcare model exhibits significant racial bias, causing Black patients to receive less medical care than others. The costs and potential risks associated with large-scale language models, such as discriminatory biases, are discussed in [5] and the research community is encouraged to consider the impacts of the ever-increasing size of DNNs beyond just the model’s accuracy for a target task.

In this paper, we showed another major concern regarding DNNs that enables attacks on the users’ privacy; even when users might think it is safe to only release a very narrow result of their private data. An important concern on the potential misuse of ML models is that unlike the discovery of software vulnerabilities that can be quickly patched, it is very difficult to propose effective defenses against harmful consequences of ML models [62]. It is suggested [28] that tech regulators should become “proactive”, rather than being “reactive”, and design controlled, confidential, and automated experiments with black-box access to ML services. While service providers may argue that ML models are proprietary resources, it is not unreasonable to allow appropriate regulators to have controlled and black-box access for regulatory purposes. For instance, a regulator can check the curiosity of ML classifiers provided by cloud APIs via a test set that includes samples with sensitive attributes.

6.2 Attacks in Machine Learning

ML models can leak detailed sensitive information about their training datasets in both white-box and black-box views [65]. Since DNNs tend to learn as many features as they can, and some of these features are inherently useful in inferring more than one attribute, then DNNs trained for seemingly non-sensitive attributes can implicitly learn other potentially sensitive attributes [15, 66]. While information-theoretical approaches [46, 51, 78] are proposed for training DNNs such that they do not leak sensitive attributes, it is shown that the empirical implementation of these theoretical approaches cannot effectively eliminate this risk [66]. A technique based on transfer learning is proposed in [66] to “re-purpose” a classifier trained for target attribute into a model for classifying a different attribute. However, the re-purposing of a classifier is different from building an HBC classifier as the former does not aim to look honest and the privacy violation is due to the further use of a classifier for other purposes without the consent of the training data owner.

In model extraction (*a.k.a.* model stealing) attack [73], an adversary aims to build a copy of a black-box ML model, without having any prior knowledge about the model’s parameters or training data and just by having access to the soft predictions provided by the model. This attack is evaluated by two objectives [25]: *test accuracy*, which measures the correctness of predictions made by the stolen model, and *fidelity*, which measures the similarity in predictions (even if it is wrong) between the stolen and the original model. Model extraction has shown the richness of a classifier’s output in reconstructing the classifier itself, but in our work we show how this richness can be used for encoding sensitive attributes at inference time.

ML enables unprecedented applications, *e.g.*, automated medical diagnosis [12, 19]. As personal data, *e.g.*, medical images, are highly sensitive, privacy-preserving learning, such as federated learning [42], is proposed to train DNN classifiers on distributed private data [39, 74]. Users who own sensitive data usually participate in training a DNN for a specified target task. Although differential privacy [11] can protect a model from memorizing its training data [1, 48], the threat model introduced in this paper is different from the commonly studied setting in property [43], or membership [59, 63] inference attacks where classifier is trained on a dataset including multiple users and the server is curious about inferring a sensitive property about users, or the presence or absence of a target user in the input dataset. We consider a threat to the privacy of a single user at inference time when the server observes only outputs of a pre-trained model.

6.3 Defense Challenges

Our experiments and analyses on the performance and behavior of HBC classifiers imply challenges in defending against such a privacy threat. First, we observed that distinguishing HBC models from standard ones is not trivial, and typical users mostly do not have the technical and computational power and resources to examine the ML services before using them. One can suggest adding random noise to the model’s outputs before sharing them [37], but, because of corresponding utility losses, users cannot just simply apply such randomization to every ML model they use. Thus, we

need robust mechanisms to discover HBC models, but also entities and systems for performing such investigations. Second, for proactive investigations, we need datasets labeled with multiple attributes, which are not always possible. Good and sufficient data is usually in the possession of ML service providers who are actually the untrusted parties in our setting. Third, users’ data might include several types of sensitive attributes, and even if we aim to distinguish HBC models we may not know which attribute an HBC model is trained for. There might be unrecognized sensitive attributes included in some of our personal data; for example, it has recently been shown [3] that DNNs can be trained to predict race from chest X-rays and CT scans of patients, in a setting where clinical experts cannot.

We believe our work serves in improving the users’ awareness about such privacy threats, and invites the community to work on efficient mechanisms as well as systems for protecting users’ privacy against such an attack.

7 CONCLUSION

We introduced and systematically studied a major vulnerability in high-capacity ML classifiers that are trained by semi-trusted ML service providers. We showed that deep classifiers can secretly encode a sensitive attribute of their private input data into their public target outputs. Our results show that, even when classifier outputs are very restricted in their form, they are still rich enough to carry information about more than one attribute. We translated this problem into an information-theoretical framework and proposed empirical methods that can efficiently implement such an attack to the privacy of users. We analyzed several properties of classifier outputs and specifically showed that the entropy of the outputs can represent the curiosity of the model up to a certain extent. Furthermore, we showed that this capability can even be transferred to other classifiers that are trained using such an HBC classifier. Finally, while we showed that even soft outputs of a multi-class classifier can be exploited for encoding sensitive information, our results suggest that it is a bit safer for a user to release soft outputs than raw outputs; without damaging the utility.

Future Work. We suggest the following open directions for further exploration. First, rigorous techniques that can help in distinguishing standard and HBC classifiers are needed, which is in the same direction of research in understanding and interpreting DNN behaviors [36]. Second, a limitation of our proposed attacks is that the sensitive attribute has to be known at training time. We suggest extending the proposed methods, or designing new methods, for encoding more than one sensitive attribute in the classifier’s output, which, in general, will be more challenging, but not impossible. Third, to investigate the implication and effects of employing an HBC model in collaborative/federated learning and multi-party ML applications, where some parties might not be fully trusted. Fourth, it is of interest to understand whether a HBC classifier reveals more information about its training dataset or less compared to standard classifiers. As we can imagine a setting where a user might also be an adversary to the server, this is a challenge for servers that utilize private data for training HBC classifiers, and looking into such scenarios will be of interest.

ACKNOWLEDGMENTS

This work was funded by the European Research Council (ERC) through Starting Grant BEACON (no. 677854) and by the UK EPSRC (grant no. EP/T023600/1) within the CHIST-ERA program. Anastasia thanks JPMorgan Chase & Co for the funding received through the J.P. Morgan A.I. Research Award 2019. Views or opinions expressed herein are solely those of the authors listed. Authors thank Milad Nasr for his help in shaping the final version of the paper.

REFERENCES

- [1] Martin Abadi, Andy Chu, Ian Goodfellow, H Brendan McMahan, Ilya Mironov, Kunal Talwar, and Li Zhang. 2016. Deep Learning With Differential Privacy. In *ACM SIGSAC Conference on Computer and Communications Security (CCS)*.
- [2] Nitin Agrawal, Ali Shahin Shamsabadi, Matt J Kusner, and Adrià Gascón. 2019. QUOTIENT: two-party secure neural network training and prediction. In *ACM SIGSAC Conference on Computer and Communications Security (CCS)*.
- [3] Imon Banerjee, Ananth Reddy Bhimireddy, John L Burns, Celi, et al. 2021. Reading Race: AI Recognises Patient’s Racial Identity in Medical Images. *arXiv:2107.10356* (2021).
- [4] Donald Beaver. 1991. Perfect Privacy for Two-Party Protocols. In *DIMACS Workshop on Distributed Computing and Cryptography*, Vol. 2.
- [5] Emily M Bender, Timnit Gebru, Angelina McMillan-Major, and Shmargaret Shmitchell. 2021. On the Dangers of Stochastic Parrots: Can Language Models be Too Big. In *Conference on Fairness, Accountability, and Transparency (FAccT)*.
- [6] Yoshua Bengio, Aaron Courville, and Pascal Vincent. 2013. Representation Learning: A Review and New Perspectives. *IEEE Transactions on Pattern Analysis and Machine Intelligence* 35, 8 (2013).
- [7] Battista Biggio, Blaine Nelson, and Pavel Laskov. 2012. Poisoning Attacks against Support Vector Machines. In *International Conference on Machine Learning (ICML)*.
- [8] Rich Caruana. 1997. Multitask Learning. *Machine Learning* 28, 1 (1997).
- [9] Benny Chor and Eyal Kushilevitz. 1991. A Zero-One Law for Boolean Privacy. *SIAM Journal on Discrete Mathematics* 4, 1 (1991).
- [10] Emiliano De Cristofaro. 2020. An Overview of Privacy in Machine Learning. *arXiv:2005.08679* (2020).
- [11] Cynthia Dwork, Aaron Roth, et al. 2014. The Algorithmic Foundations of Differential Privacy. *Foundations and Trends in Theoretical Computer Science* 9, 3-4 (2014).
- [12] Andre Esteva, Brett Kopley, Roberto A Novoa, Justin Ko, Susan M Swetter, Helen M Blau, and Sebastian Thrun. 2017. Dermatologist-Level Classification of Skin Cancer with Deep Neural Networks. *Springer Nature* 542, 7639 (2017).
- [13] European Union’s Horizon 2020 Research and Innovation Programme. 2021. Shaping the Ethical Dimensions of Smart Information Systems. <https://www.project-ssherpa.eu/>. (2021). Accessed: 2021-07-01.
- [14] Chelsea Finn, Pieter Abbeel, and Sergey Levine. 2017. Model-Agnostic Meta-Learning for Fast Adaptation of Deep Networks. In *International Conference on Machine Learning (ICML)*.
- [15] Matt Fredrikson, Somesh Jha, and Thomas Ristenpart. 2015. Model Inversion Attacks that Exploit Confidence Information and Basic Countermeasures. In *ACM SIGSAC Conference on Computer and Communications Security (CCS)*.
- [16] Craig Gentry, Amit Sahai, and Brent Waters. 2013. Homomorphic Encryption from Learning with Errors: Conceptually-Simpler, Asymptotically-Faster, Attribute-Based. In *Springer Annual Cryptology Conference*.
- [17] Oded Goldreich. 2009. *Foundations of Cryptography: Volume 2, Basic Applications*. Cambridge University Press.
- [18] Yves Grandvalet and Yoshua Bengio. 2005. Semi-supervised Learning by Entropy Minimization. In *Advances in Neural Information Processing Systems (NIPS)*.
- [19] Varun Gulshan, Lily Peng, Marc Coram, Martin C Stumpe, Derek Wu, Arunachalam Narayanaswamy, Subhashini Venugopalan, Kasumi Widner, Tom Madams, Jorge Cuadros, et al. 2016. Development and Validation of a Deep Learning Algorithm for Detection of Diabetic Retinopathy in Retinal Fundus Photographs. *Journal of American Medical Association* 316, 22 (2016).
- [20] Saurabh Gupta, Judy Hoffman, and Jitendra Malik. 2016. Cross Modal Distillation for Supervision Transfer. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*.
- [21] Song Han, Huizi Mao, and William J Dally. 2016. Deep Compression: Compressing Deep Neural Networks with Pruning, Trained Quantization and Huffman Coding. In *International Conference on Learning Representations (ICLR)*.
- [22] Geoffrey Hinton, Oriol Vinyals, and Jeffrey Dean. 2015. Distilling the Knowledge in a Neural Network. In *NIPS Workshop on Deep Learning and Representation Learning*.
- [23] Matthew Hutson. 2021. Who Should Stop Unethical A.I.? *The New Yorker Annals of Technology* (2021).
- [24] Ahmet Iscen, Giorgos Tolias, Yannis Avrithis, and Ondrej Chum. 2019. Label Propagation for Deep Semi-supervised Learning. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*.
- [25] Matthew Jagielski, Nicholas Carlini, David Berthelot, Alex Kurakin, and Nicolas Papernot. 2020. High Accuracy and High Fidelity Extraction of Neural Networks. In *USENIX Security Symposium*.
- [26] Matthew Jagielski, Alina Oprea, Battista Biggio, Chang Liu, Cristina Nita-Rotaru, and Bo Li. 2018. Manipulating Machine Learning: Poisoning Attacks and Countermeasures for Regression Learning. In *IEEE Symposium on Security and Privacy (S&P)*.
- [27] Michael Kearns and Aaron Roth. 2019. *The Ethical Algorithm: The Science of Socially Aware Algorithm Design*. Oxford University Press.
- [28] Michael Kearns and Aaron Roth. 2020. Ethical Algorithm Design Should Guide Technology Regulation. *Brookings Institution’s Artificial Intelligence and Emerging Technology Initiative* (2020).
- [29] Diederik P Kingma and Jimmy Ba. 2014. Adam: A Method for Stochastic Optimization. In *International Conference on Learning Representations (ICLR)*.
- [30] Eyal Kushilevitz, Silvio Micali, and Rafail Ostrovsky. 1994. Reducibility and Completeness in Multi-Party Private Computations. In *IEEE Annual Symposium on Foundations of Computer Science (FOCS)*.
- [31] Yann LeCun, Yoshua Bengio, and Geoffrey Hinton. 2015. Deep Learning. *Springer Nature* 521, 7553 (2015).
- [32] Tyler Lee and Anthony Ndirango. 2019. Generalization in Multitask Deep Neural Classifiers: A Statistical Physics Approach. In *Advances in Neural Information Processing Systems (NeurIPS)*.
- [33] Lemonade Insurance Company. 2021. Lemonade’s Claim Automation. <https://www.lemonade.com/blog/lemonades-sclaim-sautomation>. (2021). Accessed: 2021-07-01.
- [34] Yugeng Liu, Rui Wen, Xinlei He, Ahmed Salem, Zhikun Zhang, Michael Backes, Emiliano De Cristofaro, Mario Fritz, and Yang Zhang. 2021. ML-Doctor: Holistic Risk Assessment of Inference Attacks Against Machine Learning Models. *arXiv:2102.02551* (2021).
- [35] Ziwei Liu, Ping Luo, Xiaogang Wang, and Xiaoou Tang. 2015. Deep Learning Face Attributes in the Wild. In *International Conference on Computer Vision (ICCV)*.
- [36] Scott M Lundberg and Su-In Lee. 2017. A Unified Approach to Interpreting Model Predictions. In *Advances in Neural Information Processing Systems (NIPS)*.
- [37] Xinjian Luo, Yuncheng Wu, Xiaokui Xiao, and Beng Chin Ooi. 2021. Feature Inference Attack on Model Predictions in Vertical Federated Learning. In *IEEE International Conference on Data Engineering (ICDE)*.
- [38] David J.C. MacKay. 2003. *Information Theory, Inference and Learning Algorithms*. Cambridge University Press.
- [39] Mohammad Malekzadeh, Burak Hasircioglu, Nitish Mital, Kunal Katarya, Mehmet Emre Ozfatura, and Deniz Gündüz. 2021. Dopamine: Differentially Private Federated Learning on Medical Data. *2nd AAAI Workshop on Privacy-Preserving Artificial Intelligence (PPAI-21)* (2021).
- [40] Charles T Marx, Richard Lanus Phillips, Sorelle A Friedler, Carlos Scheidegger, and Suresh Venkatasubramanian. 2019. Disentangling Influence: Using Disentangled Representations to Audit Model Predictions. In *Advances in Neural Information Processing Systems (NeurIPS)*.
- [41] David McAllester and Karl Stratos. 2020. Formal Limitations on the Measurement of Mutual Information. In *Conference on Artificial Intelligence and Statistics (AISTATS)*.
- [42] Brendan McMahan, Eider Moore, Daniel Ramage, Seth Hampson, and Blaise Agüera y Arcas. 2017. Communication-Efficient Learning of Deep Networks from Decentralized Data. In *Conference on Artificial Intelligence and Statistics (AISTATS)*.
- [43] Luca Melis, Congzheng Song, Emiliano De Cristofaro, and Vitaly Shmatikov. 2019. Exploiting Unintended Feature Leakage in Collaborative Learning. In *IEEE Symposium on Security and Privacy (S&P)*.
- [44] Fatemehsadat Mirshghallah, Mohammadkazem Taram, Praneeth Vepakomma, Abhishek Singh, Ramesh Raskar, and Hadi Esmaeilzadeh. 2020. Privacy in Deep Learning: A survey. *arXiv:2004.12254* (2020).
- [45] Fan Mo, Ali Shahin Shamsabadi, Kleomenis Katevas, Soteris Demetriou, Ilias Leontiadis, Andrea Cavallaro, and Hamed Haddadi. 2020. DarkNet: Towards Model Privacy at the Edge using Trusted Execution Environments. In *International Conference on Mobile Systems, Applications, and Services (MobiSys)*.
- [46] Daniel Moyer, Shuyang Gao, Rob Brekelmans, Greg Ver Steeg, and Aram Galstyan. 2018. Invariant Representations without Adversarial Training. In *Advances in Neural Information Processing Systems (NeurIPS)*.
- [47] Kevin P Murphy. 2021. *Probabilistic Machine Learning: An Introduction*. MIT Press.
- [48] Milad Nasr, Reza Shokri, and Amir Houmansadr. 2018. Machine Learning with Membership Privacy Using Adversarial Regularization. In *ACM SIGSAC Conference on Computer and Communications Security (CCS)*.
- [49] Ziad Obermeyer, Brian Powers, Christine Vogeli, and Sendhil Mullainathan. 2019. Dissecting Racial Bias in an Algorithm Used to Manage the Health of Populations. *Science* 366, 6464 (2019).

- [50] Tae-Hyun Oh, Tali Dekel, Changil Kim, Inbar Mosseri, William T Freeman, Michael Rubinstein, and Wojciech Matusik. 2019. Speech2Face: Learning the Face Behind a Voice. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*.
- [51] Seyed Ali Osia, Ali Shahin Shamsabadi, Sina Sajadmanesh, Ali Taheri, Kleomenis Katevas, Hamid R Rabiee, Nicholas D Lane, and Hamed Haddadi. 2020. A Hybrid Deep Learning Architecture for Privacy-Preserving Mobile Analytics. *IEEE Internet of Things Journal* 7, 5 (2020).
- [52] Nicolas Papernot, Martin Abadi, Ulfar Erlingsson, Ian Goodfellow, and Kunal Talwar. 2016. Semi-supervised Knowledge Transfer for Deep Learning from Private Training Data. In *International Conference on Learning Representations (ICLR)*.
- [53] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas Kopf, Edward Yang, Zachary DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy, Benoit Steiner, Lu Fang, Junjie Bai, and Soumith Chintala. 2019. PyTorch: An Imperative Style, High-Performance Deep Learning Library. In *Advances in Neural Information Processing Systems (NeurIPS)*.
- [54] Andrew Pavder, Andrew Martin, and Ian Brown. 2014. Modelling and Automatically Analysing Privacy Properties for Honest-but-Curious Adversaries. *University of Oxford Technical Report* (2014).
- [55] Ben Poole, Sherjil Ozair, Aaron Van Den Oord, Alex Alemi, and George Tucker. 2019. On Variational Bounds of Mutual Information. In *International Conference on Machine Learning (ICML)*.
- [56] PyTorch. Accessed: 2021-05-01. Pruning Library. <https://github.com/pytorch/pytorch/blob/master/torch/nn/utils/prune.py>. (Accessed: 2021-05-01).
- [57] Adityanarayanan Radhakrishnan, Mikhail Belkin, and Caroline Uhler. 2020. Overparameterized Neural Networks Implement Associative Memory. *Proceedings of the National Academy of Sciences* 117, 44 (2020).
- [58] Scott E. Reed, Honglak Lee, Dragomir Anguelov, Christian Szegedy, Dumitru Erhan, and Andrew Rabinovich. 2015. Training Deep Neural Networks on Noisy Labels with Bootstrapping. In *International Conference on Learning Representations (ICLR)*.
- [59] Ahmed Salem, Yang Zhang, Mathias Humbert, Pascal Berrang, Mario Fritz, and Michael Backes. 2018. MI-Leaks: Model and Data Independent Membership Inference Attacks and Defenses on Machine Learning Models. In *Network and Distributed Systems Security Symposium (NDSS)*.
- [60] Andrew M Saxe, Yamini Bansal, Joel Dapello, Madhu Advani, Artemy Kolchinsky, Brendan D Tracey, and David D Cox. 2019. On the Information Bottleneck Theory of Deep Learning. *Journal of Statistical Mechanics: Theory and Experiment* 2019, 12 (2019).
- [61] Jürgen Schmidhuber. 2015. Deep Learning in Neural Networks: An Overview. *Elsevier Neural Networks* 61 (2015).
- [62] Toby Shevlane and Allan Dafoe. 2020. The Offense-Defense Balance of Scientific Knowledge: Does Publishing AI Research Reduce Misuse?. In *AAAI/ACM Conference on AI, Ethics, and Society*.
- [63] Reza Shokri, Marco Stronati, Congzheng Song, and Vitaly Shmatikov. 2017. Membership Inference Attacks Against Machine Learning Models. In *IEEE Symposium on Security and Privacy (S&P)*.
- [64] Ravid Shwartz-Ziv and Naftali Tishby. 2017. Opening the Black Box of Deep Neural Networks via Information. *arXiv:1703.00810* (2017).
- [65] Congzheng Song, Thomas Ristenpart, and Vitaly Shmatikov. 2017. Machine Learning Models That Remember Too Much. In *ACM SIGSAC Conference on Computer and Communications Security (CCS)*.
- [66] Congzheng Song and Vitaly Shmatikov. 2020. Overlearning Reveals Sensitive Attributes. In *International Conference on Learning Representations (ICLR)*.
- [67] Christian Szegedy, Wojciech Zaremba, Ilya Sutskever, Joan Bruna, Dumitru Erhan, Ian Goodfellow, and Rob Fergus. 2014. Intriguing Properties of Neural Networks. In *International Conference on Learning Representations (ICLR)*.
- [68] Antti Tarvainen and Harri Valpola. 2017. Mean Teachers are Better Role Models: Weight-averaged Consistency Targets Improve Semi-supervised Deep Learning Results. In *Advances in Neural Information Processing Systems (NIPS)*.
- [69] Surat Teerapittayanon, Bradley McDanel, and Hsiang-Tsung Kung. 2017. Distributed Deep Neural Networks over the Cloud, the Edge and End Devices. In *IEEE International Conference on Distributed Computing Systems (ICDCS)*. IEEE.
- [70] Tijmen Schep. 2021. How Normal Am I? <https://www.hownormalami.eu>. (2021). Accessed: 2021-07-01.
- [71] Naftali Tishby, Fernando C. Pereira, and William Bialek. 2000. The Information Bottleneck Method. In *Annual Allerton Conference on Communication, Control and Computing*.
- [72] Naftali Tishby and Noga Zaslavsky. 2015. Deep Learning and the Information Bottleneck Principle. In *IEEE Information Theory Workshop (ITW)*.
- [73] Florian Tramèr, Fan Zhang, Ari Juels, Michael K Reiter, and Thomas Ristenpart. 2016. Stealing Machine Learning Models Via Prediction APIs. In *USENIX Security Symposium*.
- [74] Stacey Truex, Nathalie Baracaldo, Ali Anwar, Thomas Steinke, Heiko Ludwig, Rui Zhang, and Yi Zhou. 2019. A Hybrid Approach to Privacy-Preserving Federated Learning. In *12th ACM Workshop on Artificial Intelligence and Security*. 1–11.
- [75] Matias Vera, Pablo Piantanida, and Leonardo Rey Vega. 2018. The Role of The Information Bottleneck in Representation Learning. In *IEEE International Symposium on Information Theory (ISIT)*.
- [76] Cunrui Wang, Qingling Zhang, Wanquan Liu, Yu Liu, and Lixin Miao. 2019. Expression of Concern: Facial Feature Discovery for Ethnicity Recognition. *Wiley Interdisciplinary Reviews: Data Mining and Knowledge Discovery* 9, 1 (2019).
- [77] Ji Wang, Weidong Bao, Lichao Sun, Xiaomin Zhu, Bokai Cao, and S Yu Philip. 2019. Private Model Compression via Knowledge Distillation. In *AAAI Conference on Artificial Intelligence*.
- [78] Ji Wang, Jianguo Zhang, Weidong Bao, Xiaomin Zhu, Bokai Cao, and Philip S Yu. 2018. Not Just Privacy: Improving Performance of Private Deep Learning in Mobile Cloud. In *ACM SIGKDD International Conference on Knowledge Discovery & Data Mining (KDD)*.
- [79] Yilun Wang and Michal Kosinski. 2018. Deep Neural Networks are More Accurate than Humans at Detecting Sexual Orientation from Facial Images. *Journal of personality and social psychology* 114, 2 (2018).
- [80] Yandong Wen, Bhiksha Raj, and Rita Singh. 2019. Face Reconstruction from Voice using Generative Adversarial Networks. In *Advances in Neural Information Processing Systems (NeurIPS)*.
- [81] Alan F Westin. 1968. Privacy and Freedom. *Washington and Lee Law Review* 25, 1 (1968).
- [82] Jacob Whitehill, Gwen Littlewort, Ian Fasel, Marian Bartlett, and Javier Movellan. 2009. Toward Practical Smile Detection. *IEEE Transactions on Pattern Analysis and Machine Intelligence* 31, 11 (2009).
- [83] Ashia Wilson. 2018. *Lyapunov Arguments in Pptimization*. Ph.D. Dissertation. UC Berkeley.
- [84] Qizhe Xie, Minh-Thang Luong, Eduard Hovy, and Quoc V Le. 2020. Self-training with Noisy Student Improves ImageNet Classification. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*.
- [85] Chao Zhang, Dacheng Tao, Tao Hu, and Bingchen Liu. 2020. Generalization Bounds of Multitask Learning From Perspective of Vector-Valued Function Learning. *IEEE Transactions on Neural Networks and Learning Systems* (2020).
- [86] Zhifei Zhang, Yang Song, and Hairong Qi. 2017. Age Progression/Regression by Conditional Adversarial Autoencoder. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*.

APPENDIX

A THE CONVEX USE CASE

In this section, we first present an example of an HBC classifier with a convex loss function, then we analyze the behavior of such a classifier.

A.1 An Example of Convex Loss

Let $\mathbf{x} = [x_0, x_1, \dots, x_{M-1}, 1]$ denote the data sample with target attribute $y \in \{0, 1\}$, and sensitive attribute $s \in \{0, 1\}$, and for each $i \in [M]$ we have $x_i \in \mathbb{R}$. Let us consider a binary logistic regression classifier \mathcal{F} that is parameterized by $\Theta \in \mathbb{R}^{M+1}$:

$$\hat{y} = \mathcal{F}(\mathbf{x}) = \sigma(\Theta^\top \mathbf{x}) = \frac{1}{1 + \exp(-\Theta^\top \mathbf{x})}.$$

To make \mathcal{F} a honest-but-curious classifier, the server can choose multipliers $\beta^y \in [0, 1]$ and $\beta^s \in [0, 1]$ such that $\beta^y + \beta^s = 1$ and build the loss function

$$\begin{aligned} \mathcal{L}^{\mathcal{F}} &= \beta^y \mathcal{L}^y + \beta^s \mathcal{L}^s \\ &= -\beta^y (y \log \hat{y} + (1-y) \log(1-\hat{y})) \\ &\quad - \beta^s (s \log \hat{y} + (1-s) \log(1-\hat{y})) \\ &= -(\beta^y y + \beta^s s) \log \hat{y} \\ &\quad - (\beta^y (1-y) + \beta^s (1-s)) \log(1-\hat{y}) \\ &= -(\beta^y y + \beta^s s) \log \hat{y} \\ &\quad - ((\beta^y + \beta^s) - (\beta^y y + \beta^s s)) \log(1-\hat{y}) \\ &= -(\beta^y y + \beta^s s) \log \hat{y} \\ &\quad - (1 - (\beta^y y + \beta^s s)) \log(1-\hat{y}) \\ &= -z \log \hat{y} - (1-z) \log(1-\hat{y}) \end{aligned} \quad (12)$$

where $z = \beta^y y + \beta^s s$. We know that:

- (1) If function \mathcal{A} is convex and $\alpha \geq 0$, then $\alpha \mathcal{A}$ is convex.
- (2) If two functions \mathcal{A}_1 and \mathcal{A}_2 are convex, then $\mathcal{A}_1 + \mathcal{A}_2$ is convex.
- (3) Functions $-\log \hat{y} = -\log(\sigma(\Theta^\top \mathbf{x}))$ and $-\log(1-\hat{y}) = -\log(1-\sigma(\Theta^\top \mathbf{x}))$, are convex w.r.t. the parameters of a logistic regression model [47].

Thus, considering these three facts, $\mathcal{L}^{\mathcal{F}}$ in Eq (12) is convex.

A.2 Analysis of Convex Loss

Now, we analyze the situations in which loss functions for training HBC classifiers \mathcal{F} , parameterized by Θ , are convex. Let the loss function for these situations be

$$\begin{aligned} \mathcal{L}^{\mathcal{F}} &= \beta^y \mathcal{L}^y(\mathcal{F}(\mathbf{x}; \Theta), y) + \beta^s \mathcal{L}^s(\mathcal{G}(\mathcal{F}(\mathbf{x}; \Theta)), s) \\ &= \beta^y \mathcal{L}^y(\Theta) + \beta^s \mathcal{L}^s(\Theta), \end{aligned}$$

where we assume \mathcal{G} is a fixed attack, and in the second line we just simplified the equation to focus on the trainable parameters Θ . Let Θ_*^y and Θ_*^s denote the optimal parameters that minimize \mathcal{L}^y and \mathcal{L}^s , respectively. Assume that, at each iteration t of training, loss functions \mathcal{L}^y and \mathcal{L}^s are μ^y - and μ^s -strongly convex functions with respect to Θ_t , respectively. Considering the gradient descent dynamics in continuous time [83],

$$d\Theta_t = -\nabla_{\Theta} \mathcal{L}^{\mathcal{F}} dt,$$

we can define the Lyapunov function $V_t = \frac{1}{2} \|\Theta_t - \Theta_*^y\|_2^2$. Observe that by strong convexity of \mathcal{L}^y , the optimum of this Lyapunov function is Θ_*^y and it is unique. Now, we can analyze the dynamics of this Lyapunov function to better understand the convergence of the gradient descent dynamics when the loss function is set to be optimizing both for y and s .

First, let us consider $\beta^s = 0$, then

$$\begin{aligned} dV_t &= (\Theta_t - \Theta_*^y)^\top d\Theta_t = -(\Theta_t - \Theta_*^y)^\top \nabla_{\Theta} \mathcal{L}^{\mathcal{F}} dt \\ &= -\beta^y (\Theta_t - \Theta_*^y)^\top \nabla_{\Theta} \mathcal{L}^y(\Theta_t) dt \\ &= -\beta^y (\Theta_t - \Theta_*^y)^\top (\nabla_{\Theta} \mathcal{L}^y(\Theta_t) - \nabla_{\Theta} \mathcal{L}^y(\Theta_*^y)) dt \\ &\leq -\beta^y \mu^y \|\Theta_t - \Theta_*^y\|_2^2 dt, \end{aligned} \quad (13)$$

where in the second line we use the gradient descent dynamics, in the third line we use the fact that $\nabla_{\Theta} \mathcal{L}^y(\Theta_*^y) = 0$ (remember that we assume Θ_*^y is the set of parameters that is optimal for \mathcal{L}^y), and in the last line we use the strong convexity of the function. By integrating both sides of inequality in (13) we obtain

$$V_t \leq e^{-t\beta^y\mu^y} V_0, \quad (14)$$

so that as t increases, we converge closer to the optimum Θ_*^y ; that is a well-known result [83].

Now, we show how the addition of the loss term on the sensitive attribute, \mathcal{L}^s , prevents the convex classifier from fully converging. Let $\beta^s > 0$, then

$$\begin{aligned} dV_t &= -(\Theta_t - \Theta_*^y)^\top \nabla_{\Theta} \beta^y \mathcal{L}^y(\Theta_t) dt - \beta^s (\Theta_t - \Theta_*^y)^\top \nabla_{\Theta} \mathcal{L}^s(\Theta_t) dt \\ &\leq -\beta^y \mu^y \|\Theta_t - \Theta_*^y\|_2^2 dt + \beta^s (\mathcal{L}^s(\Theta_*^y) - \mathcal{L}^s(\Theta_t)) dt, \end{aligned}$$

where we use convexity of \mathcal{L}^y and that by convexity³ of \mathcal{L}^s it holds $(\Theta_*^y - \Theta_t)^\top \nabla_{\Theta} \mathcal{L}^s(\Theta_t) \leq \mathcal{L}^s(\Theta_*^y) - \mathcal{L}^s(\Theta_t)$. Integrating, we obtain

$$V_t \leq e^{-t\beta^y\mu^y} V_0 + \beta^s \int_0^t e^{\beta^s \mathcal{L}^y(\Theta_t)(u-t)} (\mathcal{L}^s(\Theta_*^y) - \mathcal{L}^s(\Theta_u)) du.$$

This result shows that exact convergence is only obtained if $\Theta_*^y = \Theta_*^s$ since in this case, by the convexity of \mathcal{L}^s , the right-most term would have been upper-bounded by zero as Θ_*^y would have been the optimizer also for \mathcal{L}^s . Otherwise, we only converge to a neighborhood of the optimum Θ_*^y , where the size of the neighborhood is governed by β^s and the properties of \mathcal{L}^s .

As a summary, the trade-off that we face for a classifier with *limited* capacity is that: if we also optimize for the sensitive attribute, we will only ever converge to a neighborhood of the optimum for the target attribute. While these results hold for a convex setting and are simplistic in nature, they give us a basic intuition into the idea that when the attributes y and s are somehow *correlated*, the output \hat{y} can better encode both tasks. This result also gives us the intuition that when we have *uncorrelated* attributes, we should need classifiers with more capacity to cover the payoff for not converging to the optimal point of target attribute.

B PSEUDOCODE OF PARAMETERIZED ATTACK

In Algorithm 1 we show the pseudo-code of the training process we presented in Section 4.3.1.

³In a convex function, for all x and y we have: $f(x) \geq f(y) + f'(y)(x - y)$.

Algorithm 1 Training an HBC classifier and its corresponding attack based on information bottleneck formulation in Section 4.3.1

```

1: Input:  $\mathcal{F}$ : the model chosen as the classifier,  $\mathcal{G}$ : the model chosen as the attack,  $\mathbb{D}^{train}$ : training dataset,  $(\beta^x, \beta^y, \beta^s)$ : trade-off multipliers,  $E$ : number of training epochs,  $K$ : batch size.

2: Output: Updated  $\mathcal{F}$  and  $\mathcal{G}$ .
3: Random initialization of  $\mathcal{F}$  and  $\mathcal{G}$ .
4: for  $e : 1, \dots, E$  do
5:   for  $b : 1, \dots, |\mathbb{D}^{train}|/K$  do
6:      $(X, Y, S) \leftarrow$  a random batch of  $K$  samples  $(x, y, s) \sim \mathbb{D}^{train}$ 

7:      $\hat{Y} = \mathcal{F}(X)$ 
8:      $\hat{S} = \mathcal{G}(\hat{Y})$ 
9:      $\hat{H}(S|\hat{Y}) = -\sum_{k=1}^K \sum_{i=0}^{S-1} \mathbb{I}_{(s^k=i)} \log(\hat{s}_i^k)$ 
10:    Update  $\mathcal{G}$  via the gradients of loss function  $\mathcal{L}^{\mathcal{G}} = \hat{H}(S|\hat{Y})$ 
11:     $\hat{S} = \mathcal{G}(\hat{Y})$ 
12:     $\hat{H}(S|\hat{Y}) = -\sum_{k=1}^K \sum_{i=0}^{S-1} \mathbb{I}_{(s^k=i)} \log(\hat{s}_i^k)$ 
13:     $\hat{H}(Y|\hat{Y}) = -\sum_{k=1}^K \sum_{i=0}^{Y-1} \mathbb{I}_{(y^k=i)} \log(\hat{y}_i^k)$ 
14:     $\hat{H}(\hat{Y}) = -\sum_{k=1}^K \sum_{i=0}^{Y-1} \hat{y}_i^k \log(\hat{y}_i^k)$ 
15:    Update  $\mathcal{F}$  via the gradients of loss function  $\mathcal{L}^{\mathcal{F}} = \beta^x \hat{H}(\hat{y}) + \beta^y \hat{H}(y|\hat{y}) + \beta^s \hat{H}(s|\hat{y})$ 
16:   end for
17: end for

```

C ADDITIONAL EXPERIMENTAL RESULTS

In Table 7 we investigate the effect of the classifier’s capacity on its performance and behavior. Results show that both honesty (δ^y) and curiosity (δ^s) decrease by reducing the capacity of classifier in a similar way. A more interesting observation is that the average of the outputs’ entropy tends to increase when reducing the classifier’s capacity, which is due to both having lower honesty and also the emerged difficulties because of curiosity. This gives us the hint that by having a higher-capacity classifier, a server cannot only achieve better honesty-curiosity trade-offs, but also can easier hide the suspicious behavior of the classifier’s output.

Table 8 and Table 9 show experimental results for settings where target and sensitive attribute are reversed, compared to Table 3 and Table 4 in Section 5. We can observe similar patterns in these results as well, confirming our findings explained in Section 5.

D CURIOSITY AND GENERALIZATION

Underlying the capabilities of the HBC classifiers is the ability of the output \hat{y} to encode multiple attributes. Such capability also forms the foundation for *representation learning* [6], *multi-task learning* [14] and *disentanglement of factors* [40]. It is well-known that added noise during training, e.g., in the form of Dropout, improves the generalization performance of a model [75]. These observations coincide with results obtained in [43], where it is shown that differential privacy is not effective protection against property inference attacks; even more so, noise may improve the performance by keeping the model from overfitting on the main task.

Information bottleneck (IB) is used for understanding the nature of intermediate representations produced by the hidden layers of

Table 7: The effect of classifier’s capacity on the UTKFace dataset for the target attribute Age and the sensitive attribute Race where $Y = S = 3$. Similar to Table 5, we also show the average of the entropy of classifier’s output by $\tilde{H}(\hat{y})$ in bits. We fix $(\beta^x, \beta^y, \beta^s) = (.0, .7, .3)$. Model capacity is determined in the first column, where a “Layer Size” of A-B-C-D-E refers to the number of neurons in each layer, respectively (related to the “Output Size” in Table 12). We show results in two setting: (A) RawHBC and (B) SoftHBC.

Layers’ Size	δ^y	δ^s	$\tilde{H}(\hat{y})$
(A) RawHBC			
128-128-64-64-128	81.17±0.27	83.14±0.48	0.61±0.03
64-64-32-32-64	80.33±0.20	81.52±0.12	0.73±0.03
42-42-21-21-42	78.45±0.79	79.50±0.38	0.86±0.03
32-32-16-16-32	78.05±0.62	77.82±0.33	0.86±0.03
25-25-12-12-25	75.73±0.61	74.96±0.80	1.00±0.11
21-21-10-10-21	74.59±0.97	74.01±0.60	0.98±0.11
18-18-9-9-18	73.94±0.63	71.39±0.72	1.12±0.02
16-16-8-8-16	70.80±0.96	69.75±0.90	1.14±0.06
12-12-6-6-12	64.68±5.02	66.39±2.65	1.31±0.05
10-10-5-5-10	62.17±2.51	63.77±0.87	1.40±0.05
8-8-4-4-8	62.74±2.97	63.20±1.20	1.41±0.03
(B) SoftHBC			
128-128-64-64-128	80.66±0.28	73.08±0.63	0.73±0.02
64-64-32-32-64	79.96±0.32	72.67±0.20	0.80±0.03
42-42-21-21-42	79.07±0.40	71.42±0.63	0.83±0.03
32-32-16-16-32	77.29±0.43	69.87±1.04	0.88±0.02
25-25-12-12-25	75.95±0.13	68.30±1.06	0.95±0.02
21-21-10-10-21	75.00±0.68	67.41±1.04	1.00±0.04
18-18-9-9-18	74.22±0.70	66.71±1.25	1.08±0.06
16-16-8-8-16	73.53±0.42	67.89±0.87	1.21±0.05
12-12-6-6-12	70.55±0.48	65.98±0.66	1.26±0.06
10-10-5-5-10	66.46±2.59	63.24±1.17	1.30±0.05
8-8-4-4-8	63.97±1.49	61.72±2.35	1.36±0.05

DNNs during training [60, 64, 72]. In Section 4.3.1, we discussed that formulation in Eq (5) shows a trade-off: compressing the information included in x , while encoding as much information about y and s as possible into \hat{y} . In other words, IB states that the optimal HBC classifier must keep the information in x that is useful for y and s while compressing other, irrelevant information. On the other hand, a useful classifier is one that does not overfit on its training dataset, meaning that the *generalization gap*⁴ is small.

One can notice that under sufficient “similarity” between the target and sensitive attribute, the addition of the sensitive attribute into the classifier’s optimization objective can even *improve* the generalization capabilities for the classification of target attribute [32]. Specifically, the addition of another attribute s while training for the target attribute y can reduce the generalization gap as the addition

⁴The difference between the classifier’s performance on its training data and the classifier’s performance on unseen test data is called the generalization gap of classifier.

Table 8: The honesty (δ^y) and curiosity (δ^s) of an HBC classifier trained via regularized attack, where target attribute (y) is Race and sensitive attribute (s) is Age. Values are in percentage (%).

		S = 2			S = 3		S = 4		S = 5	
		δ^y	δ^s		δ^y	δ^s	δ^y	δ^s	δ^y	δ^s
Y = 2	NC	(1.,0)	87.63 ± .19	62.37 ± .41	87.63 ± .19	51.47 ± .28	87.63 ± .19	37.80 ± .56	87.63 ± .19	32.26 ± .07
	RawHBC	(.7,.3)	87.70 ± .26	81.42 ± .86	87.82 ± .19	79.58 ± .84	88.17 ± .12	65.67 ± .79	87.94 ± .18	58.15 ± 1.6
	SoftHBC	(.7,.3)	87.22 ± .19	67.22 ± .40	87.57 ± .33	54.03 ± .78	87.26 ± .16	40.86 ± .22	87.24 ± .34	36.60 ± 2.5
Y = 3	NC	(1.,0)	85.42 ± .46	62.25 ± .87	85.42 ± .46	51.97 ± .43	85.42 ± .46	38.09 ± .77	85.42 ± .46	32.34 ± .45
	RawHBC	(.7,.3)	86.11 ± .42	82.14 ± .35	86.13 ± .31	78.09 ± .88	86.33 ± .08	65.56 ± .94	86.11 ± .34	58.95 ± 1.3
	SoftHBC	(.7,.3)	85.25 ± .26	72.47 ± 2.4	85.36 ± .61	64.52 ± 1.0	85.31 ± .38	53.03 ± .90	85.19 ± .10	45.77 ± .51
Y = 4	NC	(1.,0)	81.57 ± .33	64.99 ± .70	81.57 ± .33	54.12 ± .58	81.57 ± .33	41.13 ± .15	81.57 ± .33	35.56 ± .47
	RawHBC	(.7,.3)	81.69 ± .12	81.33 ± .38	81.08 ± .57	77.57 ± .54	81.26 ± .35	64.44 ± .54	81.15 ± .22	58.93 ± .82
	SoftHBC	(.7,.3)	81.20 ± .21	75.17 ± .35	81.27 ± .48	69.14 ± .07	81.09 ± .10	55.63 ± .11	81.21 ± .08	51.28 ± .52
Y = 5	NC	(1.,0)	80.66 ± .24	66.29 ± .39	80.66 ± .24	56.64 ± .25	80.66 ± .24	44.42 ± .76	80.66 ± .24	38.13 ± .73
	RawHBC	(.7,.3)	81.05 ± .40	80.93 ± .78	80.78 ± .66	78.48 ± .26	80.00 ± .45	65.16 ± .46	80.87 ± .70	58.71 ± .24
	SoftHBC	(.7,.3)	80.94 ± .35	79.59 ± .27	80.59 ± .08	76.57 ± .28	80.37 ± .09	62.76 ± .47	80.52 ± .06	55.75 ± .52
		(.5,.5)	80.03 ± .22	81.46 ± .23	72.36 ± .55	78.57 ± .17	70.48 ± 1.6	66.32 ± .41	70.88 ± .66	58.02 ± 1.2

Table 9: Gender and Race. Classification accuracy (%) on UTKFace test set for Race vs. Gender. In each experiment, there are 4 different tasks based on the size Y and S . Empty cells (–) are due to the incapability of regularized attacks for $S > 2$.

		S = 2		S = 3		S = 4		S = 5			
Setting	Attack	(β^y, β^s)	δ^y	δ^s	δ^y	δ^s	δ^y	δ^s	δ^y	δ^s	
Std	raw [66]	(1., 0.)	90.23 ± .21	58.62 ± 1.2	90.23 ± .21	47.54 ± 1.1	90.23 ± .21	44.21 ± .43	90.23 ± .21	43.94 ± .04	
Y = 2	RawHBC	Parameterized	(.7, .3)	90.18 ± .11	86.91 ± .32	90.56 ± .23	81.66 ± .24	90.24 ± .09	74.04 ± .48	90.28 ± .35	74.01 ± .89
		Parameterized	(.7, .3)	90.59 ± .07	56.21 ± .00	89.92 ± .11	46.26 ± .50	90.37 ± .19	44.31 ± .21	90.14 ± .34	44.93 ± .94
		Parameterized	(.5, .5)	89.74 ± .43	61.60 ± 4.1	86.95 ± .66	52.85 ± 2.3	87.69 ± .48	50.57 ± .53	87.57 ± .76	50.88 ± .78
SoftHBC		(.3, .7)	79.96 ± .80	81.54 ± .60	72.41 ± .55	77.28 ± .45	72.97 ± 7.3	63.90 ± 2.6	72.97 ± 7.3	63.90 ± 2.6	
		(.7, .3)	90.18 ± .17	73.20 ± .10	–	–	–	–	–	–	
		Regularized	(.5, .5)	89.71 ± .06	82.28 ± .10	–	–	–	–	–	
		(.3, .7)	87.86 ± .40	84.35 ± .48	–	–	–	–	–	–	

of another attribute can act as a *regularizer*, which places an inductive bias on the learning of the target attribute and guides the model towards learning more general and discerning features [8, 85].

Previous work (e.g., Theorem 1 in [75]) has shown that the generalization gap can be upper-bounded by the amount of compression happening in the classifier, i.e., the mutual information between the internal representations learned by the classifier and the input data. Basically, the compression forces the classifier to throw away information unrelated to the target task (e.g., noises). Moreover, [32] shows that training a classifier for more than one attribute can bring a major benefit in situations where the single-attribute classifier “underperforms” on the target task, e.g., due to the lack of samples in the training data or the presence of noise. Particularly this happens when attributes are related—where the relatedness is measured through the alignment of the input features—and the data for the second attribute is of good quality [32].

Our observation from Section 5 is that, while properties Age, Race, and Gender are not necessarily correlated properties, they are all coarse-grained properties for which understanding of the full image is required (as opposed to fine-grained properties present only in part of the image). As observed from the results, in certain

settings, including curiosity in the classifier can improve the performance of the main task. Therefore, because all properties are sufficiently coarse-grained, learning in a multi-task framework can keep the network from overfitting details by enforcing it to learn certain general, coarse-grained features and thus perform better in the main task.

Putting all together, it can be expected that by training an HBC classifier, we can not only reduce the generalization gap for the target attribute, but also compress the learned representations, which means decrease (or at least not increase) the entropy of the representations. Thus, in situations where the curiosity of the classifier *improves* generalization, the recognition of an HBC classifier is even more challenging.

E MORE MOTIVATIONAL EXAMPLES

As some other motivational examples, we can consider:

- (1) A sentiment analysis application that can be run on the user’s browser and, taking a text, it can help the user to estimate how their text may sound to the readers: e.g., $\hat{y} = [p_{\text{positive}} \cdot p_{\text{neutral}} \cdot p_{\text{negative}}]$.



Figure 10: Sample images from CelebA [35]. Columns (from left to right) show Male, Smiling, Attractive, High Cheekbones, and Heavy Makeup attributes. The first and second rows show Black Hair with and without the other corresponding column’s attribute, respectively. Similarly, the third and fourth rows Blond Hair, and fifth and sixth rows Brown Hair attributes.

The server may ask for collecting only outputs, *e.g.*, to improve their service, but server might try to infer if the text includes specific sensitive words or if the text falls into a certain category; by making an HBC sentiment classifier.

(2) An activity recognition app, from an insurance company, can be run on the user’s smartwatch and analyzing motion data. This app can provide an output $\hat{y} = [p_{walk}, p_{run}, p_{sit}, p_{sleep}]$. The server may ask for only collecting \hat{y} to offer the user further health statistics or advice, but it may also try to infer if the user’s body mass index (BMI) is more than or below a threshold.

(3) Automatic speech recognition (ASR) is the core technology for all voice assistants, such as Amazon Alexa, Apple Siri, Google Assistant, Microsoft Cortana, to provide specialized services to their users, *e.g.*, playing a song, reporting weather condition, or ordering food. ASR models, at the top, have a softmax layer that estimates a probability distribution on a vocabulary with a particular length, *e.g.*, 26 letters in the English alphabet plus some punctuation symbols. Service providers usually perform post-processing on the output of the ASR model, that is the probable letter for each time frame, to combine them and predict the most probable word and further the actual sentence uttered by the user. Despite the fact that they can be asked to run their ASR model at the user’s side, but

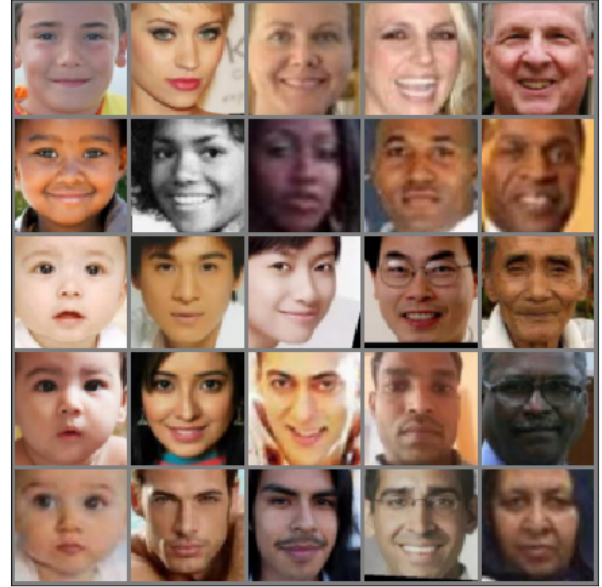


Figure 11: Sample images from UTKFace [86]. Each row shows images labeled with the same Race (White, Black, Asian, Indian, and others), and each column shows one of the five Age groups ([0-19, 20-26, 27-34, 35-49, 50-100]).

they can argue that the post-processing ASR’s outputs is crucial for providing a good service, and these outputs are supposed to only predict what the user is asking for (*i.e.*, the uttered sentence) and not any other information, *e.g.*, the Age, Gender, or Race of the user. However, an HBC ASR might secretly violate its user’s privacy.

(4) An EU-funded project [13] designed an app [70] for people to experience how ML models judge their faces. This app asked for access to the user’s camera to capture their face image. Then, pre-trained ML classifiers are being run on the user’s own computer, in the browser. The project promises users that “no personal data is collected” but at the end, users can voluntarily “share some anonymized data” [70]. Despite the useful purpose of this app, we show that even such anonymous collection of the outputs of pre-trained classifiers can result in some privacy leakage.

F DATASETS

F.1 Details of CelebA Dataset

Figure 10 shows some sample images from CelebA dataset [35]. While CelebA does not have any attributes with a number of classes more than 2, we utilize three mutually exclusive binary attributes BlackHair, BlondHair, BrownHair to build a 3-class attribute of HairColor. This still gives us a dataset of more than 115K images that fall into one of these categories.

F.2 Details of UTKFace Dataset

Figure 11 shows some sample images from UTKFace dataset [86]. Each image, collected from the Internet, has three attributes: (1) Gender (male or female), (2) Race (White, Black, Asian, Indian, or others),

Table 12: The implemented model as classifier \mathcal{F} .

Layer Type	Output Size	Number of Parameters
Conv2d(kernel=2, stride=2)	128, 32, 32	1,664
BatchNorm2d	128, 32, 32	256
LeakyReLU(slope=0.01)	128, 32, 32	0
Dropout(p=0.2)	128, 32, 32	0
Conv2d(kernel=2, stride=2)	128, 16, 16	65,664
BatchNorm2d	128, 16, 16	256
LeakyReLU(slope=0.01)	128, 16, 16	0
Dropout(p=0.2)	128, 16, 16	0
Conv2d(kernel=2, stride=2)	64, 8, 8	32,832
BatchNorm2d	64, 8, 8	128
LeakyReLU(slope=0.01)	64, 8, 8	0
Dropout(p=0.2)	64, 8, 8	0
Conv2d(kernel=2, stride=2)	64, 4, 4	16,448
BatchNorm2d	64, 4, 4	128
LeakyReLU(slope=0.01)	64, 4, 4	0
Dropout(p=0.2)	64, 4, 4	0
Linear	128	131,200
LeakyReLU(slope=0.01)	128	0
Dropout(p=0.5)	128	0
Linear	Y	128×Y+Y

Table 13: The implemented model as attack \mathcal{G} .

Layer Type	Output Size	Number of Parameters
Linear	20×Y	80×Y
LeakyReLU(slope=0.01)	20×Y	0
Dropout(p=0.25)	20×Y	0
Linear	10×Y	610×Y
LeakyReLU(slope=0.01)	10×Y	0
Dropout(p=0.25)	10×Y	0
Linear	S	10×Y×S+S

Table 10: The details of how we assign Age label to each image in UTKFace.

Label	Number of Classes				Samples
	2	3	4	5	
0	a ≤ 30	a ≤ 20	a ≤ 21	a ≤ 19	4593
1	30 > a	20 < a ≤ 35	21 < a ≤ 29	19 < a ≤ 26	5241
2		35 < a	29 < a ≤ 45	26 < a ≤ 34	4393
3			45 < a	34 < a ≤ 49	4491
4				49 < a	4987

Table 11: The details of how we assign Race label to each image in UTKFace. W:White, B:Black, A:Asian, I:Indian, or O:others.

Label	Number of Classes				Samples
	2	3	4	5	
0	W	W	W	W	10078
1	BAIO	A	B	B	4526
2		BIO	I	A	3434
3			AO	I	3975
4				O	1692

and (3) Age (from 0 to 116). Gender is a binary label and there are 12391 vs. 11314 images with male vs. female labels. Table 10 and Table 11 show how we assign Age and Race labels to each image in UTKFace for different number of classes used in our evaluations. For the Age label we choose categories such that classes become balanced, in terms of the number of samples. For the Race label, we are not that free (like Age), thus we choose categories such that classes do not become highly unbalanced and we try to keep samples in each class as similar to each other as possible.

G MODEL ARCHITECTURES

For all experiments reported in this paper, we use PyTorch [53] that is an open-source library for the implementation of deep neural networks [31, 61]. Table 12 and Table 13 show the details of neural network architecture that are implemented as classifier \mathcal{F} and attack \mathcal{G} (see Figure 6).