



Biswas, C., Ganguly, D., Mukherjee, P., Bhattacharya, U. and Hou, Y. (2022) Privacy-aware supervised classification: An informative subspace based multi-objective approach. *Pattern Recognition*, 122, 108301. (doi: [10.1016/j.patcog.2021.108301](https://doi.org/10.1016/j.patcog.2021.108301)).

This is the Author Accepted Manuscript.

There may be differences between this version and the published version. You are advised to consult the publisher's version if you wish to cite from it.

<http://eprints.gla.ac.uk/264601/>

Deposited on: 03 February 2022

Enlighten – Research publications by members of the University of Glasgow  
<http://eprints.gla.ac.uk>

# Privacy-Aware Supervised Classification: An Informative Subspace based Multi-Objective Approach

Chandan Biswas <sup>a,\*</sup>, Debasis Ganguly<sup>b</sup>, Partha Sarathi Mukherjee<sup>c</sup>,  
Ujjwal Bhattacharya<sup>a</sup>, Yufang Hou<sup>d</sup>

<sup>a</sup>*Indian Statistical Institute, Kolkata, India*

<sup>b</sup>*University of Glasgow, UK*

<sup>c</sup>*Tatras Data, Delhi, India*

<sup>d</sup>*IBM Research Europe, Dublin, Ireland*

---

## Abstract

Sharing the raw or an abstract representation of a labelled dataset on cloud platforms can potentially expose sensitive information of the data to an adversary, e.g., in the case of an emotion classification task from text, an adversary-agnostic abstract representation of the text data may eventually lead an adversary to identify the demographics of the authors, such as their gender and age. In this paper, we propose a universal defense mechanism against such malicious attempts of stealing sensitive information from data shared on cloud platforms. More specifically, our proposed method employs an informative subspace based multi-objective approach to obtain a sensitive information aware encoding of the data representation. A number of experiments conducted on both standard text and image datasets demonstrate that our proposed approach is able to reduce the effectiveness of the adversarial

---

\*Corresponding author

*Email address:* `chandanbiswas08_r@isical.ac.in` (Chandan Biswas )

task (i.e., in other words is able to better protect the sensitive information of the data) without significantly reducing the effectiveness of the primary task itself.

*Keywords:* Privacy preserving representation learning, Informative subspace, Multi-objective learning, Defence against information stealing adversarial attacks

---

## 1 Introduction

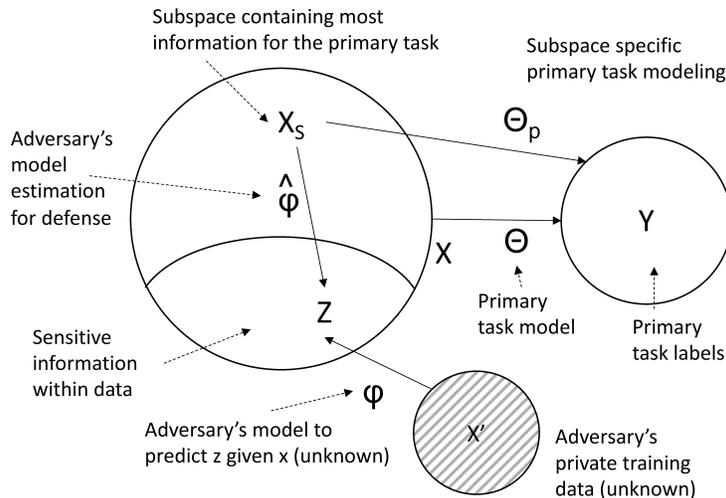
The era of data-driven learning is continuously witnessing increased computational requirements for training multi-layered complex neural networks for supervised machine learning (ML) through a layered approach of abstraction from the raw data, e.g., the work on contextual word vectors pre-trained on large collections of documents to capture the inherent language model in text [1], or that of training deep image networks to capture higher levels of visual features from images [2].

One standard solution to mitigate the intensive computational requirements of training data-driven models is to follow the standard ‘software as a service’ paradigm, in which the computations to train an ML model are provided as a service (MLaaS) by a powerful computing device (server), virtually accessible through a distributed computing environment (cloud) [3]. An MLaaS-based solution requires a user (client program) to upload an encoded form of the data, usually corresponding to an abstract representation of it, e.g. pre-trained vectors such as BERT [1] for text, or Inception-Net for

17 images [2]), to the server. Although such an MLaaS based workflow allows  
18 provision for distributed data sharing and also reduces the computational  
19 overhead of the client workstations, a risk with an MLaaS architecture is  
20 that it can potentially lead to breaches in data security and privacy [4].

21 To illustrate the point on potential threats on data privacy, consider an  
22 *adversarial model* which is able to eavesdrop on the communication channel  
23 between a client and the server offering computation on encoded forms of  
24 data. Imagine a situation where an adversarial model is *pre-trained* on past  
25 data, which in terms of its domain and characteristics, is similar to the one  
26 that is transmitted to the server over a communication channel. In such a  
27 situation, this pre-trained adversarial model could use this submitted data  
28 as an input to predict a number of *sensitive* attribute values from this data  
29 [5].

30 As a concrete example of an adversarial attack on data privacy, con-  
31 sider that the encoded data sent from a client workstation to a computation  
32 server over a communication channel corresponds to that of movie reviews,  
33 and the *primary task* for which the computational resources of the server is  
34 sought, refers to the task of classifying a review into positive or negative,  
35 i.e. the primary task involves learning a mapping of the form  $\theta : \mathbf{x} \mapsto y$ ,  
36  $\mathbf{x} \in \mathbb{R}^d, y \in \{0, 1\}$ , where  $\mathbf{x}$  represents an encoding of the data, e.g. a se-  
37 quence encoding of the words comprising the review [6]. Imagine that each  
38 review contains additional *identity information attributes*,  $z$ , corresponding  
39 to sensitive information about the author, e.g. the age, gender etc. De-



**Figure 1:** Schematic of our proposed defence mechanism that relies on identifying a candidate subspace,  $X_s$ , of the input space, on which the set of primary task labels,  $Y$ , is likely to exhibit a strong functional dependence. The remaining subspace,  $X - X_s$ , is then useful to estimate a likely functional dependence with the sensitive information,  $\hat{\phi}$ , an inversion on which is then used to defend against an adversarial model,  $\phi$ .

40 spite not being a part of the encoding, the adversary can potentially feed  
 41 the encoded data as input into an adversarial network, that has already been  
 42 trained on pairs of movie reviews encoding and the attribute values (e.g. gen-  
 43 der),  $(\mathbf{x}', z)$ , to learn an association between the two of the form  $\phi : \mathbf{x}' \mapsto z$ ,  
 44  $\mathbf{x}' \in \mathbb{R}^d, z \in \{0, 1\}$ . The parameters of the trained network,  $\phi$ , may then  
 45 accurately predict the demographics of the current encoded data  $\mathbf{x}$ , i.e., the  
 46 closer  $\mathbf{x}$  is to  $\mathbf{x}'$  the higher is the associated risk of leaking the attribute value  
 47 information [7].

48 A standard approach to prevent an attacker stealing the sensitive infor-  
 49 mation from data is to make the encoding process itself aware of the inten-  
 50 tions of an adversary, which usually involves first formulating the adversarial  
 51 model,  $\phi : \mathbf{x} \mapsto z$  as a secondary task, and then applying a multi-objective

52 based encoding transformation of the data, where the first objective corre-  
53 sponds to the primary task and the subsequent ones correspond to one or  
54 more secondary tasks, each such secondary task representing an adversarial  
55 objective [5]. The learning objective, in this case, seeks to minimize the po-  
56 tential degradation of the primary task effectiveness due to the noise which is  
57 required to be incorporated within the data as a defence against adversarial  
58 attacks.

59 **Our Contributions.** We now enlist our contributions in this paper. First,  
60 contrary to a standard approach of data-driven encoding that uses uniform  
61 weights for the abstract features, we hypothesize that the defence mecha-  
62 nism of a multi-objective based approach can potentially be improved by  
63 a weighted distribution over features. Specifically, this involves leveraging  
64 information from *candidate subspaces*,  $\mathbf{x}_s \in \mathbb{R}^k$ , ( $k < d$ ) of the input data  
65 that are *strongly correlated with the primary category labels* in the form  
66  $\theta_p : \mathbf{x}_s \mapsto y$ . The *residual subspace* is thus likely to be functionally asso-  
67 ciated to the latent attribute values of the data, or in other words, to the  
68 secondary (adversarial) task categories  $\hat{\phi} : \mathbf{x}_s' \mapsto z$ ,  $\mathbf{x}_s' \in \mathbb{R}^{d-k}$ , which in turn  
69 approximately models the function  $\phi : \mathbf{x} \mapsto z$ . We argue that this way of  
70 modeling the adversarial information yields a more robust encoding mecha-  
71 nism that is likely to be more resilient to security threats and our experiments  
72 confirm this hypothesis.

73 Second, in contrast to most existing approaches which conduct experi-  
74 ments mostly on text data with annotated metadata information (such as

75 the demographic attributes, e.g., age and gender annotated as a part of the  
76 TrustPilot dataset [5]), we report empirical results on both images and text.  
77 For images, we test our method both on implicit and explicit demographic  
78 attributes. As implicit attributes, we use stylistic attributes, such as the  
79 slant or ligatures in handwriting, that could potentially reveal the age of a  
80 person. As explicit attributes, we test if the metadata information of age  
81 and gender associated with a set of lesion images can potentially be revealed  
82 to information stealing attacks.

## 83 2. Related Work

84 **Adversarial Learning.** An adversarial attack broadly refers to the meth-  
85 ods of generating samples (often called adversarial examples) that are in-  
86 distinguishable from samples drawn from the true data distribution with an  
87 objective to ‘fool’ a classifier [8]. These attacks typically use first order gradi-  
88 ent information, such as FGSM [8], I-FGSM [9], MI-FGSM [10], Ada-FGSM  
89 [11] etc. Successful demonstrations of black-box adversarial perturbations at-  
90 tacks leading to degrading the effectiveness of classifiers were demonstrated  
91 in [12] and [13]. Defence mechanisms against such adversarial attacks include  
92 those of using regularized FGSM [14], and defensive distillation [15].

93 Different from adversarial learning, we rather employ a multi-objective  
94 encoding, the purpose of which is to ensure that it potentially would be  
95 difficult for an adversary to use a pre-trained system (on similar data) to  
96 effectively predict the values of sensitive attributes (e.g., age, gender etc.)

97 from the encoded data.

98 **Differential Privacy and Privacy-preserving Data Encoding.** The  
99 objective of differential privacy is somewhat similar to that of privacy-preserving  
100 encoding. However, differential privacy does not involve encoding the raw  
101 data as vectors; instead, it obfuscates parts of relational data so as to mit-  
102 igate individual data leakage [16]. Various de-identification or anonymizing  
103 technologies have been proposed to protect data privacy, which often involve  
104 adding noise or masking sensitive information in the released dataset [17].  
105 The concept of additive noise in differential privacy for relational databases  
106 also finds applications in Bayesian risk minimization in general [18], or in  
107 Bayesian linear regression [19] in particular. Privacy preserved data encod-  
108 ing finds applications in encoding raw data for both unsupervised [20] and  
109 supervised learning tasks [21]. For text data, privacy-preserving based en-  
110 coding is particularly crucial because the inherent characteristics of natural  
111 language (e.g., writing style or word usage patterns) often reveal information  
112 about the authors, which can be used by adversaries to reveal such sensitive  
113 information. As examples, the authors of [22] used online behavior, stylistic  
114 choices and language models to predict the age group of blog authors, while  
115 those of [23] used Twitter content to predict the occupational class.

116 A number of recent studies has proposed the dual objective of privacy  
117 preservation (minimizing leakage of sensitive information) and model preser-  
118 vation (maximizing the performance of an algorithm on the encoded data),  
119 e.g., applying a ‘multi-detasking’ model to train an adversarial classifier

120 simultaneously with the primary downstream text classifier, where during  
121 training, the primary classifier updates its parameters to confuse the attacker  
122 model [5]. The study reported in [24] developed a distributed framework for  
123 privacy preserving multi-task learning protocol by applying encryption mech-  
124 anisms. The authors of [4] explored an adversarial learning approach that  
125 learns unbiased representations of text with respect to specific sensitive at-  
126 tributes. Somewhat different from the findings of [5], the authors of [25]  
127 showed that despite adversarial training methods being generally effective in  
128 reducing the amount of implicit sensitive information, in some cases, how-  
129 ever, a substantial amount of sensitive information still persists and can be  
130 extracted from the encoded representations.

131 Although our proposed method falls into the general class of multi-objective  
132 approaches, such as those of [5] and [26], our proposed method is more general  
133 in the sense that we leverage the candidate subspaces that are most informa-  
134 tive of the primary task. Since parts of these subspaces are less likely to be  
135 comprised of the sensitive information in data, our method seeks to address  
136 some of the concerns pointed out in [25], i.e. removal of sensitive attributes  
137 (e.g. demographics) from data instances can still lead to an adversary pre-  
138 dicting this missing information. Our subspace-based approach is explicitly  
139 directed towards mitigating this problem in the sense that the privacy-aware  
140 encoding process puts more emphasis only on those components of the data  
141 that are more useful for the primary task, while suppressing the residual  
142 space that contains most of the information on the sensitive attributes.

143 **Feature Importance for Explanations.** Standard approaches of model-  
144 agnostic instance-wise explanations for classification include those of em-  
145 ploying linear regression to learn a simplified decision boundary by sampling  
146 points around a data instance [27], applying a Gumbel distribution to esti-  
147 mate instance-wise feature importance [28] etc. The authors of [29] reiterate  
148 the importance of feature selection for supervised learning tasks, whereas  
149 those of [30] and [31] explore feature selection for the case of unsupervised  
150 learning.

151 In the context of our work, we use the idea of exploring informative  
152 candidate subspaces with a parameterized approach, as first proposed in [28].  
153 An explicit use of feature importance also provides an interpretable way of  
154 preserving data privacy.

### 155 **3. A General Framework for Privacy-Aware Encoding**

156 In this section, we formally describe a general framework for defence  
157 against adversarial threats using a multi-task learning based workflow. We  
158 present a general approach to the problem in the sense that the overall frame-  
159 work allows provision to incorporate more than one adversarial task, each  
160 corresponding to a particular attribute of the data.

#### 161 *3.1. Privacy-Agnostic Encoding*

162 Using the notations introduced Section 1, the predictive model for the  
163 primary task, generally speaking, can be *learned* with a set of linear trans-

164 formation functions (realized with a multi-layer perceptron) of the form

$$\begin{aligned}
 P(y = i | \mathbf{w}; \theta, \theta_p) &= \sigma(\theta_p \cdot \mathbf{x})_i = \frac{\exp(\theta_{p_i} \cdot \theta \cdot \mathbf{w})}{\sum_{j=1}^c \exp(\theta_{p_j} \cdot \theta \cdot \mathbf{w})}, \\
 \mathbf{x} &= \theta \cdot \mathbf{w}, \mathbf{x} \in \mathbb{R}^s, \mathbf{w} \in \mathbb{R}^d, y \in \mathbb{Z}_c,
 \end{aligned}
 \tag{1}$$

165 where  $\mathbf{w} \in \mathbb{R}^d$  denotes a  $d$ -dimensional vector representation (encoding) of  
 166 the input data,  $y \in \mathbb{Z}_c$  denotes a class label (one of  $c$  possible values) corre-  
 167 sponding to the classification task,  $\theta \in \mathbb{R}^{s \times d}$  denotes a matrix of parameters  
 168 (a latent layer of a neural network), and  $\theta_p \in \mathbb{R}^{c \times s}$  denotes a matrix of pa-  
 169 rameters specifically corresponding to the classification task ( $\theta_{p_i} \in \mathbb{R}^s$  is the  
 170 parameter vector for the  $i$ -th class). As a simplification, we do not explic-  
 171 itly include the bias parameter as a part of the softmax equations. Since  
 172 the encoding process of Equation 1 does not explicitly take account an ad-  
 173 versarial threat against a subset of data attributes, the encoding  $\mathbf{x} \in \mathbb{R}^s$  is  
 174 privacy-agnostic.

### 175 3.2. Privacy-Aware Encoding

176 An encoding space different from Equation 1 that explicitly addresses a  
 177 set of sensitive attributes has been shown to be effective in defence against  
 178 adversarial models [5]. However, the work in [5] addresses the defence mech-  
 179 anism for a single attribute only. Instead, we present a more general setup  
 180 involving more than one attribute.

181 In the context of our work, the attributes manifest themselves as an im-  
 182 plicit part of the data, or otherwise, it is straight-forward to remove the

183 attributes before encoding the data [25]. In particular, we assume that the  
 184 encoding of an input data instance,  $\mathbf{w}$ , is a function of both the raw data it-  
 185 self, (say  $w$ ) and its latent characteristics (sensitive attributes). We represent  
 186 a pair comprising an input data instance and a set of  $M$  sensitive attributes  
 187 (assuming categorical values) associated with it as  $(w, \{z_1, \dots, z_M\})$ , where  
 188  $z_m \in \mathbb{Z}_{s_m}$ , i.e. there are a total of  $s_j$  number of possible values for the  $j^{th}$   
 189 attribute.

190 A multi-objective transformation then uses the pairs,  $(w, \{z_1, \dots, z_M\})$ , to  
 191 encode the privacy-agnostic representation  $\mathbf{w} \in \mathbb{R}^d$  as learnable parameters,  
 192  $\mathbf{x} \in \mathbb{R}^s$ , with the combined objective

$$\begin{aligned}
 P(y = i, z_1, \dots, z_M | \mathbf{w}; \theta, \theta_p, \phi^1, \dots, \phi^M) = \\
 (1 - \sum_{m=1}^M \gamma_m) \sigma(\theta_p \cdot \mathbf{x})_i - \sum_{m=1}^M \gamma_m \sigma(\phi^m \cdot \mathbf{x})_{z_m},
 \end{aligned} \tag{2}$$

193 where  $\mathbf{x} = \theta \cdot \mathbf{w}$ ,  $\mathbf{x} \in \mathbb{R}^s$ , and  $\mathbf{w} \in \mathbb{R}^d$ , and similar to Equation 1,  $\sigma(\cdot)_i$   
 194 is an abbreviation for the softmax function with respect to the  $i$ -th class.  
 195 The multi-objective loss of Equation 2 can be realized with a feed-forward  
 196 network comprising a shared layer (parameter matrix  $\theta \in \mathbb{R}^{s \times d}$ ) and the task  
 197 specific layers. Separate layers, one for each adversarial task ( $\phi^m \in \mathbb{R}^{s_m \times s}$ ),  
 198 in addition to the primary task itself ( $\theta_p \in \mathbb{R}^{c \times s}$ ), are all connected to the  
 199 shared layer. Note that the parameters corresponding to  $\mathbf{w}$ 's in Equation 2  
 200 are obtained from pre-trained representations and hence are not learnable.

201 To illustrate Equation 2 with an example, consider a text classification

202 problem, where each document is associated with the demographic attributes  
203 - age ( $z_1$ ) and gender ( $z_2$ ) of author. In such a situation, the value of  $M$  in  
204 Equation 2 would be 2. Continuing with the example, if age is discretized  
205 into 3 categories, e.g., ‘young’, ‘middle-aged’ and ‘senior’ then  $s_1 = 3$ .

206 In a generalized setting, the multi-objective loss function of Equation 2  
207 models a relative trade-off between the effectiveness of the primary task and  
208 the desired lack of effectiveness of the adversarial ones (notice the negative  
209 factor in the linear combination corresponding to the adversarial tasks). A  
210 low value of each linear combination parameter,  $\gamma_m \in [0, 1] : (\sum_m \gamma_m <$   
211  $1)$ , associates a small importance to the necessity of defending against an  
212 information stealing attack against the  $m$ -th attribute. Notice that setting  
213  $\gamma_m = 0$  degenerates Equation 2 to the privacy-agnostic encoding of Equation  
214 1.

#### 215 4. An Information Theoretic Perspective

216 In this section, we describe how to extend the general multi-task based  
217 privacy preserving approach from an information theoretic perspective. As  
218 per the motivation behind the schematic depiction of Figure 1, we now for-  
219 mally describe how to leverage information from the importance of features  
220 (components of the encoded vector representation of a data instance) to help  
221 the process of learning a better encoding for privacy preservation.

222 *4.1. Subspace Encoding*

223 A limitation of Equation 2 is that the parameters of the shared layer and  
224 the primary-task specific layer (i.e.  $\theta$  and  $\theta_p$  respectively) are trained with  
225 respect to the entire feature space of the encoded vector  $\mathbf{w}$ , whereas it is  
226 more likely to be the case that a part of this feature space correlates strongly  
227 with the primary task. The key idea in our proposed method is to substitute  
228 the encoding  $\mathbf{w}$  of Equation 2 with a subset of features that are most likely  
229 to be informative for the primary task. This has a two-fold advantage.

230 First, a subspace of the most informative features for the primary task is  
231 likely to lead to a down-weighting of the residual subspace potentially con-  
232 stituting information responsible for determining the values of the sensitive  
233 attributes of the data. In other words, this is likely to degrade the effective-  
234 ness of the secondary tasks thus providing a potentially improved defence  
235 mechanism.

236 Second, since the subspace-based encoding approach puts more emphasis  
237 on parts of the data that are potentially responsible for determining the  
238 primary task output, it is also likely to lead to improving the effectiveness of  
239 the primary task itself.

240 *4.2. Parameterized Subspace Selection with Gumbel Distribution*

241 The authors of [28] computed the importance of features by measuring  
242 the mutual information between the primary task labels and an arbitrary  
243 feature subspace  $\mathbf{w}_s \in \mathbb{R}^k$ , ( $k < d$ ). The total number of possible subspaces,

244  $\binom{d}{k}$ , is exponential for relatively large values of  $k$ . Hence finding an opti-  
 245 mal subspace representing the largest amount of information for data driven  
 246 models is a challenging problem. A solution, proposed in [28, 32], is to use  
 247 a parameterized version of a subspace (specifically obtained with a Gumbel  
 248 distribution) that allows a gradient descent based optimization of its param-  
 249 eters. The objective is seek an optimum state of maximum informativeness  
 250 of the subspace with respect to a set of labels. Before describing how this  
 251 is applied in the context of our problem, we present a brief overview of the  
 252 Gumbel based learning of subspaces, mostly following the exposition of [28].

253 A Gumbel distribution,  $G(0, 1)$ , is a distribution of random variables of  
 254 the form  $G_i = -\log(-\log u_i)$ ,  $u_i \sim \mathcal{U}(0, 1)$ ,  $\mathcal{U}$  being the uniform distribution.  
 255 The Gumbel softmax probability distribution uses a concrete distribution,  
 256 which is a continuous differentiable approximation of a categorical random  
 257 variable. The *Gumbel softmax* is a modification of the softmax function  
 258 involving random variables sampled from the Gumbel distribution, one each  
 259 for each component of the softmax. In the context of our problem, we use the  
 260 Gumbel softmax distribution to estimate the importance of each component  
 261 of the encoding vector,  $\mathbf{w} \in \mathbb{R}^d$ . Formally speaking,

$$C = \{C_i : C_i = \frac{\exp((\log w_i + G_i)/\rho)}{\sum_{j=1}^d \exp((\log w_j + G_j)/\rho)}, i = 1, \dots, d\}, \quad (3)$$

262 where  $\rho$  is a *temperature* parameter, higher values of which makes the dis-  
 263 tribution close to uniform (for our experiments, we set  $\rho = 0.1$  as per [28]).

264 To select  $k$  features from a set of available  $d$  features, one needs to inde-  
 265 pendently sample  $k$  times from the Gumbel softmax distribution resulting in  
 266 a total of  $k$  random vectors  $\{\mathbf{c}_1, \dots, \mathbf{c}_k\}$ , where the  $j^{th}$  vector  $\mathbf{c}_j$  is sampled  
 267 from Gumbel softmax, i.e.,  $\mathbf{c}_j \sim C$ . Let  $\Lambda_k \in \mathbb{R}^{d \times k}$  be the matrix constituted  
 268 from the  $k$  random vectors,  $\mathbf{c}_j$ , thus sampled. A row-wise maximum of the  
 269 matrix,  $\Lambda_k$  then yields an approximation of a  $k$ -hot random vector  $\lambda_k \in \mathbb{R}^d$ .  
 270 The highest  $k$  elements of  $\lambda_k$  (corresponding to the most important features)  
 271 are retained while the rest ( $d - k$ ) are set to 0. Thus  $\lambda_k$  is a vector with  
 272  $k$  non-zero elements (*soft  $k$ -hot*) determining the choice of a  $k$ -dimensional  
 273 subspace.

#### 274 4.3. Feature Subspace with Multi-Objective

275 In the context of our problem (see Equation 2), data is represented as  
 276 vectors in  $d$  dimensions, i.e.  $\mathbf{w} \in \mathbb{R}^d$ , out of which we intend to select a  
 277 subspace  $\mathbf{w}_s \in \mathbb{R}^k$  comprised of the most informative features. After selecting  
 278 a random vector with  $k$  non-zero elements,  $\lambda_k$ , we now model its interaction  
 279 with the primary classification task as

$$\begin{aligned}
 P(y = i, z_1, \dots, z_M | \mathbf{w}; \theta, \theta_p, \phi^1, \dots, \phi^M) = \\
 (1 - \sum_{m=1}^M \gamma_m) \sigma(\theta_p \cdot \mathbf{x})_i - \sum_{m=1}^M \gamma_m \sigma(\phi^m \cdot \mathbf{x})_{z_m},
 \end{aligned} \tag{4}$$

280 where  $\mathbf{x} = \theta \cdot (\mathbf{w} \odot \lambda_k)$ ,  $\mathbf{x} \in \mathbb{R}^s$  and  $\mathbf{w} \in \mathbb{R}^d$ . Equation 4 is a more con-  
 281 strained form of Equation 2. This is because instead of considering an arbi-  
 282 trary  $s$ -dimensional transformation from  $\mathbf{w}$  (privacy-agnostic encoding) to  $\mathbf{x}$

283 (privacy-aware encoding) of Equation 2, we specifically select an informative  
284 subspace, denoted by, say  $\mathbf{w}_s = \mathbf{w} \odot \lambda_k$ . This is obtained by an element-wise  
285 multiplication of the input encoding with a soft  $k$ -hot vector obtained from  
286 the Gumbel softmax distribution.

287 As a next step, the informative subspace is used to learn the privacy-  
288 aware encoded representation<sup>1</sup>. In our experiments, instead of specifying the  
289 value of  $k$  directly, we control it with a fraction,  $\tau \in [0, 1]$  of the input data  
290 dimension, i.e.,  $k = \lfloor \tau d \rfloor$ .

## 291 5. Experimental Setup

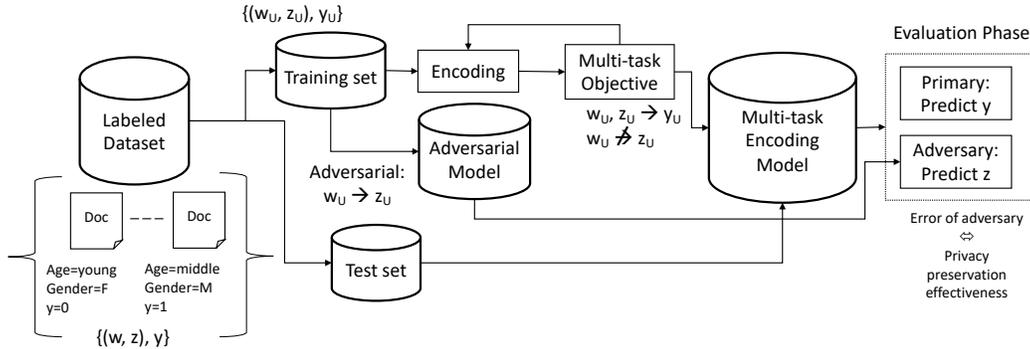
### 292 5.1. Experiment Workflow

293 A laboratory based setup is devoid of the presence of a true adversary (e.g.  
294 as shown in the schematic of Figure 1). In such a situation, the adversary  
295 would have access to a pre-trained model which is trained to predict the  
296 sensitive attributes from input data instances. An adversarial model is likely  
297 to be more harmful if it has been trained on data instances that resemble  
298 the ones (i.e. similar in terms of encoded vector representations) to the ones  
299 that are sent over from the client to the MLaaS. To mimic this situation as  
300 closely as possible in a laboratory setup, we set up our experiments as shown  
301 in Figure 2.

302 For each labeled dataset, each data instance is annotated with additional

---

<sup>1</sup>A prototype of the implementation is available for research purposes at <https://github.com/chandanbiswas08/12x-mt>



**Figure 2:** Schematics of the common setup for the evaluation workflow. Both the privacy-aware encoding and the adversarial model (one for each attribute) is trained on the train-split of the data. During evaluation phase, the *privacy-preserved* encoded vectors for the test-split are fed into the adversarial model to predict values of the attributes. The prediction error of this pseudo-adversarial setup indicates the effectiveness of privacy preservation.

303 attribute value pairs. With this we train a logistic regression model on the  
 304 train-split of the data to simulate an adversarial attack of predicting these  
 305 additional attribute values from the data (a separate adversarial model is  
 306 trained for each attribute type, shown as a single model in Figure 2 to avoid  
 307 clutter).

308 In general, corresponding to  $M$  different attribute types (see Equations  
 309 2 and 4), we evaluate the effectiveness of the adversarial task as an *inverse*  
 310 *effectiveness* measure for a particular defence method used in our experi-  
 311 ments. The experiment workflow ensures that the encoding process of a  
 312 defence mechanism is oblivious of the category values (e.g., values of age and  
 313 gender) of the test-split.

**Table 1:** Summary of the dataset used in our experiments.

Dataset	#Instances		Primary task		Adversarial Tasks		
	Train	Test	Classify	#Classes	Type	Attribute	Categories
Morpho-MNIST (M-MNIST)	120K	40K	Digits	10	Synthetic	Slant Fracture	{left, neutral, right} {yes, no}
Skin Cancer MNIST (HAM10K)	8500	1500	Diseases	7	Real	Age Gender	{ $\leq 30$ , 31-60, $> 60$ } {male, female}
Trustpilot (US English)	23K	4K	Sentiment	2	Real	Age Gender	{ $\leq 35$ , $> 35$ } {male, female}

## 314 5.2. Dataset

315 A dataset suitable for the purpose of our experiments needs to be an-  
316 notated with additional attribute values corresponding to the sensitive in-  
317 formation, the prediction of which during the adversarial workflow branch  
318 (see Figure 2) could then be set up as information leakage. To test the ef-  
319 fectiveness of our proposed subspace based privacy preservation approach  
320 on different modalities of data, we experiment with both text and image  
321 datasets. The details of each dataset follows next (also summarized in Table  
322 1).

323 **Morpho-MNIST (M-MNIST).** The primary task of the original MNIST  
324 dataset involves detecting the class of a digit (a gray-scale image with  $28 \times 28$   
325 pixels) out of the 10 possibilities (one of 0 to 9). As a part of latent infor-  
326 mation that can potentially be leaked from an encoding of a hand-written  
327 image (e.g. a 2d convolution with maxpooling), we first consider the *slant* of  
328 a hand-written digit, which can be considered to be correlated with person-  
329 ality traits [33]. To setup the dataset, each *slant label*,  $z_1$  (in our notation),  
330 is obtained by applying a threshold on the horizontal shear,  $\alpha$ . The value of  
331 the shear,  $\alpha$ , in turn is computed as a function of second order moments of

332 the gray-scale values,  $x_{ij}$  [34]. Formally,

$$z_1 = \begin{cases} 0 & \alpha \leq -0.3 \text{ (left)} \\ 1 & -0.3 < \alpha < 0.3 \text{ (neutral)} \\ 2 & \alpha \geq 0.3 \text{ (right)} \end{cases} \quad (5)$$

333 In addition to the slant, the second attribute that we address in our ex-  
334 periments is whether the image of a hand-written digit is *fractured*, i.e., a  
335 lack of continuity is exhibited in the strokes. The value of this attribute, if re-  
336 vealed in a real-life situation, could indicate the age of an OCR-ed document  
337 to an adversary.

338 For our experiments with the fracture attribute, we use an existing dataset,  
339 namely the ‘Morpho-MNIST’, where morphological erosion is applied to syn-  
340 thetically generate fractured images [34]. Addition of the synthetically gener-  
341 ated fractured images, one for each image in the original MNIST, resulted in  
342 doubling the number of images for this dataset. The information on whether  
343 an image is fractured is not available to an adversary, nor does the adversary  
344 is allowed to compute the slant labels using Equation 5.

345 **Skin Cancer MNIST (HAM10K)**. Contrary to using synthetically gen-  
346 erated attribute values for the adversarial task, the ‘Skin Cancer MNIST’  
347 (or HAM10K) dataset [35] allows us to setup the adversarial tasks with two  
348 explicitly annotated attributes. The primary task in this dataset involves  
349 identifying one out of 7 possible skin diseases, e.g., Bowen’s disease, basal



**Figure 3:** Left to right: No slant and fractures, followed by fractures with neutral, left and right slants.



**Figure 4:** Left to right: Lesion images of a young female, mid-aged male, old female and an old male.

350 cell carcinoma etc., from images of lesions. The objective in this case is to  
351 encode the data in such a way that it does not reveal the age or gender of a  
352 person without substantially degrading the effectiveness of the primary task.  
353 Some sample images from the two image datasets are shown in Figures 3 and  
354 4.

355 **TrustPilot Dataset.** For the text modality, we use the TrustPilot reviews  
356 (the US English subset). The primary task on this dataset involves identify-  
357 ing sentiment (positive or negative) of a review [36]. This dataset, comprised  
358 of over 27K reviews with sentiment score ranging between 1 and 5, has an-  
359 notated values for both age and gender. Since the number of reviews with  
360 scores 2 and 3 is substantially small, we binarize the sentiment class labels  
361 by thresholding with a value of 3, i.e. scores from 1-3 are mapped to class  
362 0 and the rest to 1. Following the previous experiment setup of [5] and  
363 [4], we binarize the attribute ‘age’ as young ( $\text{age} \leq 35$ ) and its complement  
364 (representing the category ‘not young’).

365 *5.3. Baselines*

366 As baselines, we compare the following approaches. First, we apply a  
367 privacy agnostic logistic regression based approach (see Equation 1), which  
368 we denote as **LR**. Our next baseline, denoted as **MT**, is the multi-tasking  
369 based approach from existing literature [5], which we presented in this paper  
370 as Equation 2. To explore if subspace based information usage, which forms  
371 a part of our proposed method, is indeed effective, we conduct experiments  
372 with two ablation baselines. The first of these baselines (applicable for text)  
373 involves the following. After computing the term feature weights with a  
374 simple term importance statistics (specifically tf-idf), for each sentence we  
375 retain only a fraction,  $\tau \in [0, 1]$ , of the terms with the highest weights. The  
376 rationale of this baseline, denoted as **LR-TFIDF**, is to see if removing a  
377 subset of features, not correlated to the primary task alone, can prevent  
378 information leakage of secondary attributes.

379 The second ablation baseline is a degenerate case of Equation 4, where  
380 we set  $\gamma_m = 0$  for each adversarial task. This means that the  $k$ -dimensional  
381 encoding of the data, being agnostic of the adversarial tasks, only takes into  
382 account the informative subspace of the primary task. Unlike LR-TFIDF,  
383 this baseline method, denoted as **L2X** in our experiments, is applied to both  
384 text and images.

385 *5.4. Evaluation Metrics and Parameters*

386 As an evaluation metric, we employ a combination of the primary task  
387 accuracy (higher the better) and the inverse accuracy of the secondary tasks  
388 (lower the better). A high value of the combined metric reflects a better  
389 defence against information leakage without a substantial drop in primary  
390 task effectiveness. For combination, we specifically use the harmonic mean  
391 between the inverse of the aggregated accuracy values of the secondary tasks  
392 and the accuracy of the primary task, i.e.,

$$F_S = \frac{2A_P(1 - A_S)}{(1 - A_S) + A_P}, \quad (6)$$

393 where each  $A_S$  is the harmonic mean over the accuracy of each adversarial  
394 task,  $A_{S_i}$ .

395 The hyper-parameters tuned for each method were: a)  $\tau$ , which controls  
396 the number of features retained (for the LR-TFIDF baseline, this refers to  
397 the fraction of the terms retained with the highest tf-idf scores), and b)  
398  $(\gamma_1, \gamma_2)$ , which controls the relative importance of the two adversarial tasks  
399 (Equation 4). In particular, the range of these hyper-parameters in our  
400 experiments were:  $[0.2, 0.8]$  for  $\tau$ , and  $[0.1, 0.4]$  for  $\gamma_1$  and  $\gamma_2$ , in steps of 0.2  
401 and 0.1 respectively.

402 *5.5. Results*

403 **Summary.** Table 2 summarizes the optimal results of the different privacy  
404 preservation learning methods. The optimal result for each method was

**Table 2:** Privacy preservation results summary on different datasets. Parameter combinations that are not applicable for a method are shown as filled up gray cells. e.g. the parameter  $\gamma_1$  for LR.

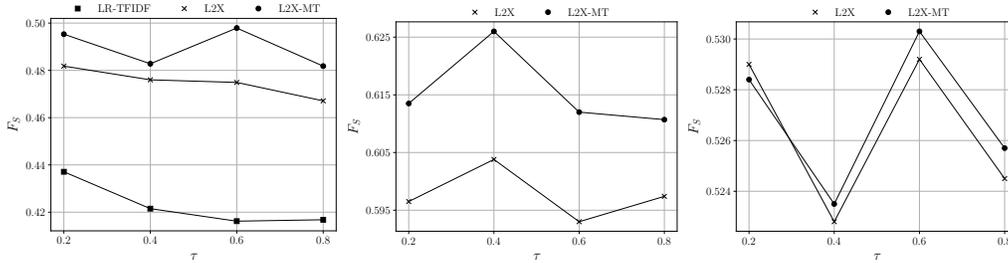
Dataset	Method	Hyper-parameters			Accuracy			Combined Measures		
		$\tau$	$\gamma_1$	$\gamma_2$	$A_P$	$A_{S_1}$	$A_{S_2}$	$F_{S_1}$	$F_{S_2}$	$F_S$
TrPilot	LR				0.8674	0.7292	0.7168	0.4127	0.4270	0.4200
	LR-TFIDF	0.2			0.8194	0.7113	0.6928	0.4270	0.4469	0.4371
	MT		0.4	0.4	0.8694	0.6849	0.6920	0.4626	0.4549	0.4587
	L2X	0.2			<b>0.8726</b>	0.6804	0.6546	0.4678	0.4949	0.4818
	L2X-MT	0.6	0.1	0.1	0.8711	<b>0.6564</b>	<b>0.6465</b>	<b>0.4928</b>	<b>0.5029</b>	<b>0.4979</b>
M-MNIST	LR				0.9840	0.8956	0.6992	0.1888	0.4608	0.3525
	MT		0.2	0.2	<b>0.9851</b>	0.8647	0.6735	0.2379	0.4904	0.3896
	L2X	0.4			0.9593	0.5435	0.5764	0.6186	0.5877	0.6038
	L2X-MT	0.4	0.4	0.1	0.9596	<b>0.5291</b>	<b>0.5420</b>	<b>0.6318</b>	<b>0.6201</b>	<b>0.6260</b>
HAM10K	LR				0.6995	0.5757	0.6256	0.5282	0.4877	0.5093
	MT		0.3	0.2	<b>0.7072</b>	0.5749	0.6249	0.5310	0.4902	0.5119
	L2X	0.2			0.6861	0.5384	0.6045	0.5519	0.5018	0.5290
	L2X-MT	0.6	0.4	0.4	0.6861	<b>0.5376</b>	<b>0.6017</b>	<b>0.5525</b>	<b>0.5040</b>	<b>0.5303</b>

405 obtained by individually tuning its hyper-parameters.

406 We observe that although LR, being a privacy agnostic approach, results  
 407 in high effectiveness for the primary task classification, it also yields high  
 408 values for the adversarial tasks. This indicates a substantial information  
 409 leakage with the LR method. Multi-tasking based encoding (MT) helps  
 410 improve results, specially for text, as also noted in [5].

411 Subspace encoding alone (L2X) is also able to decrease the accuracy val-  
 412 ues for the adversarial tasks (i.e. improve privacy preservation), which also  
 413 means that a combination of MT and L2X should also improve results. This  
 414 is precisely what is demonstrated by the results of our method (L2X-MT),  
 415 which yields the best results for each dataset.

416 **Parameter Sensitivity.** We also investigate the effects of varying  $\tau$  (sub-  
 417 space selection), and the relative importance of the adversarial task ( $\gamma_m$ ) pa-  
 418 rameters (Equations 2 and 4) on the overall effectiveness of privacy-preservation

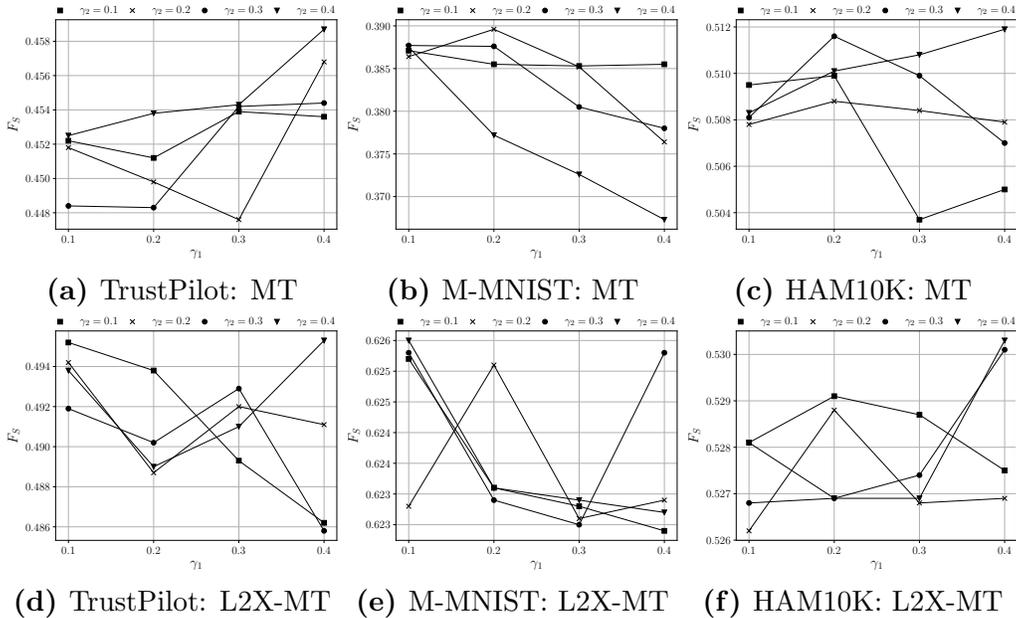


**Figure 5:** Sensitivity of the privacy-aware learning approaches with respect to relative subspace dimensionality  $\tau$ ; Left: TrustPilot, Middle: M-MNIST, Right: HAM10K.

419 learning of the corresponding primary tasks. Figure 5 shows that L2X-MT  
 420 outperforms the baselines consistently for a range of different subspace di-  
 421 mensions. Figure 6 shows the relative comparisons between the two multi-  
 422 tasking approaches - MT and L2X-MT. It can be seen that for a range of  
 423 different relative importance of the two adversarial tasks (e.g. age/gender  
 424 detection for Trustpilot and HAM10K, and slant/fracture detection for M-  
 425 MNIST), leveraging information from informative subspaces helps improve  
 426 the overall balance between primary task effectiveness and prevention of in-  
 427 formation leakage.

428 In summary, our experiments revealed the following two key observations.

- 429 1. Learning on data encoded by our method yields comparable results  
 430 with that obtained on data in its original form, i.e. *our proposed en-*  
 431 *coding does not lead to a significant decrease in the effectiveness of a*  
 432 *classification model.*
- 433 2. Data encoded by our method *significantly reduces the effectiveness of*  
 434 *an adversarial classification model* which seeks to predict sensitive at-  
 435 tributes from the data. It is also shown that the use of the *informative*



**Figure 6:** Sensitivity of MT, L2X-MT with variations in relative importance of two adversarial tasks.

436            *subspace helps to improve the defence mechanism*, i.e., it further reduces  
 437            the effectiveness of the adversarial classification model.

## 438 6. Conclusions and Future Work

439            We proposed a generic method of privacy-preserving supervised learning,  
 440            which is potentially beneficial for distributing an encoding of the input data  
 441            over a cloud environment with the end-goal of eventually learning a predictive  
 442            model (primary task) on the data. Our generic methodology combines the  
 443            advantages of two main hypotheses - that of (a) using a *multi-task objective*  
 444            that in addition to learning the primary task also learns the complemen-  
 445            tary (inverse) characteristics of an adversarial model as a defence mechanism  
 446            against information stealing attacks; and (b) using a *residual subspace* of the

447 data to further improve the defence mechanism.

448 Our experiments on image and textual data demonstrated that our pro-  
449 posed method, which jointly learns a multi-objective encoding over informa-  
450 tive subspaces (with respect to the primary task), outperforms a separate  
451 application of each.

452 In future, we would like to explore how may it be possible to obtain a  
453 privacy-preservation encoding of the input data in those cases where the sen-  
454 sitive attributes are latent rather than being manifested as explicitly anno-  
455 tated identifiable attributes (i.e., to address the situation when the attribute  
456 value annotations are not available in the training set). Unsupervised analy-  
457 sis of the input space coupled with a semi-supervised encoding approach can  
458 potentially be useful to tackle such a situation.

## 459 **References**

- 460 [1] J. Devlin, M.-W. Chang, K. Lee, K. Toutanova, BERT: Pre-training  
461 of deep bidirectional transformers for language understanding, in: Pro-  
462 ceedings of the 2019 Conference of the North American Chapter of the  
463 Association for Computational Linguistics: Human Language Technolo-  
464 gies, Volume 1 (Long and Short Papers), Association for Computational  
465 Linguistics, Minneapolis, Minnesota, 2019, pp. 4171–4186.
- 466 [2] C. Szegedy, S. Ioffe, V. Vanhoucke, A. Alemi, Inception-v4, inception-  
467 resnet and the impact of residual connections on learning, in: Proceed-  
468 ings of the AAAI Conference on Artificial Intelligence, Vol. 31, 2017.

- 469 [3] M. Ribeiro, K. Grolinger, M. A. Capretz, Mlaas: Machine learning as  
470 a service, in: 2015 IEEE 14th International Conference on Machine  
471 Learning and Applications (ICMLA), IEEE, 2015, pp. 896–902.
- 472 [4] Y. Li, T. Baldwin, T. Cohn, Towards robust and privacy-preserving  
473 text representations, in: Proceedings of the 56th Annual Meeting of the  
474 Association for Computational Linguistics (Volume 2: Short Papers),  
475 Association for Computational Linguistics, Melbourne, Australia, 2018,  
476 pp. 25–30.
- 477 [5] M. Coavoux, S. Narayan, S. B. Cohen, Privacy-preserving neural repre-  
478 sentations of text, in: Proc. of EMNLP '18, 2018, pp. 1–10.
- 479 [6] Q. Le, T. Mikolov, Distributed representations of sentences and docu-  
480 ments, in: Proc. of ICML'14, 2014, pp. II–1188–II–1196.
- 481 [7] B. Weggenmann, F. Kerschbaum, Syntf: Synthetic and differentially  
482 private term frequency vectors for privacy-preserving text mining, in:  
483 ACM SIGIR '18, 2018, pp. 305–314.
- 484 [8] I. J. Goodfellow, J. Shlens, C. Szegedy, Explaining and harnessing adver-  
485 sarial examples, in: Y. Bengio, Y. LeCun (Eds.), ICLR '15, San Diego,  
486 CA, USA, May 7-9, 2015, Conference Track Proceedings, 2015.
- 487 [9] A. Kurakin, I. J. Goodfellow, S. Bengio, Adversarial examples in the  
488 physical world, in: ICLR '17, Toulon, France, April 24-26, 2017, Work-  
489 shop Track Proceedings, 2017.

- 490 [10] Y. Dong, F. Liao, T. Pang, H. Su, J. Zhu, X. Hu, J. Li, Boosting  
491 adversarial attacks with momentum, in: Proc. of CVPR '18, 2018, pp.  
492 9185–9193.
- 493 [11] Y. Shi, Y. Han, Q. Zhang, X. Kuang, Adaptive iterative attack towards  
494 explainable adversarial robustness, Pattern Recognition (2020) 107309.
- 495 [12] N. Papernot, P. McDaniel, I. Goodfellow, S. Jha, Z. B. Celik, A. Swami,  
496 Practical black-box attacks against machine learning, in: Proceedings of  
497 the 2017 ACM on Asia Conference on Computer and Communications  
498 Security, ASIA CCS '17, Association for Computing Machinery, New  
499 York, NY, USA, 2017, p. 506519.
- 500 [13] D. Li, J. Zhang, K. Huang, Universal adversarial perturbations against  
501 object detection, Pattern Recognition 110 (2021) 107584.
- 502 [14] C. Szegedy, W. Zaremba, I. Sutskever, J. Bruna, D. Erhan, I. J. Goodfel-  
503 low, R. Fergus, Intriguing properties of neural networks, in: Y. Bengio,  
504 Y. LeCun (Eds.), Proc. of ICLR'14, 2014.
- 505 [15] N. Papernot, P. McDaniel, X. Wu, S. Jha, A. Swami, Distillation as a  
506 defense to adversarial perturbations against deep neural networks, in:  
507 2016 IEEE Symposium on Security and Privacy (SP), IEEE, 2016, pp.  
508 582–597.
- 509 [16] C. Dwork, Differential privacy, in: 33rd International Colloquium on  
510 Automata, Languages and Programming, part II (ICALP 2006), Vol.

- 511 4052 of Lecture Notes in Computer Science, Springer Verlag, 2006, pp.  
512 1–12.
- 513 [17] R. Wang, B. C. Fung, Y. Zhu, Q. Peng, Differentially private data pub-  
514 lishing for arbitrarily partitioned data, *Information Sciences* 553 (2021)  
515 247–265.
- 516 [18] C. Dimitrakakis, B. Nelson, Z. Zhang, A. Mitrokotsa, B. I. Rubinstein,  
517 Differential privacy for bayesian inference through posterior sampling,  
518 *JMLR* 18 (1) (2017) 343–381.
- 519 [19] G. Bernstein, D. R. Sheldon, Differentially private bayesian linear re-  
520 gression, in: *Proc. of NIPS '19*, 2019, pp. 525–535.
- 521 [20] C. Biswas, D. Ganguly, D. Roy, U. Bhattacharya, Privacy preserving  
522 approximate k-means clustering, in: *Proc. of CIKM '19*, 2019, pp. 1321–  
523 1330.
- 524 [21] Y. Jinfeng, W. Jun, J. Rong, Privacy and regression model preserved  
525 learning., in: *Proc. of AAAI '14*, 2014, pp. 1341–1347.
- 526 [22] S. Rosenthal, K. McKeown, Age prediction in blogs: A study of style,  
527 content, and online behavior in pre- and post-social media generations,  
528 in: *Proceedings of the 49th Annual Meeting of the Association for Com-  
529 putational Linguistics: Human Language Technologies*, 2011, pp. 763–  
530 772.

- 531 [23] D. Preoțiuc-Pietro, V. Lampos, N. Aletras, An analysis of the user oc-  
532 cupational class through twitter content, in: Proceedings of the 53rd  
533 Annual Meeting of the Association for Computational Linguistics and  
534 the 7th International Joint Conference on Natural Language Processing  
535 (Volume 1: Long Papers), Association for Computational Linguistics,  
536 Beijing, China, 2015, pp. 1754–1764.
- 537 [24] K. Liu, N. Uplavikar, W. Jiang, Y. Fu, Privacy-preserving multi-task  
538 learning, in: Proc. of ICDM '18, IEEE, 2018, pp. 1128–1133.
- 539 [25] Y. Elazar, Y. Goldberg, Adversarial removal of demographic attributes  
540 from text data, in: Proc. of EMNLP '18, 2018, pp. 11–21.
- 541 [26] P. Sen, D. Ganguly, Towards socially responsible ai: Cognitive bias-  
542 aware multi-objective learning, in: Proceedings of the AAAI Conference  
543 on Artificial Intelligence, Vol. 34, 2020, pp. 2685–2692.
- 544 [27] S. M. Lundberg, S.-I. Lee, A unified approach to interpreting model  
545 predictions, in: Proc. of NIPS '17, 2017, pp. 4765–4774.
- 546 [28] J. Chen, L. Song, M. Wainwright, M. Jordan, Learning to explain: An  
547 information-theoretic perspective on model interpretation, in: Proc. of  
548 ICML '18, 2018, pp. 883–892.
- 549 [29] S. Gao, G. Ver Steeg, A. Galstyan, Variational information maximiza-  
550 tion for feature selection, in: Proc. of NIPS '16, 2016, pp. 487–495.

- 551 [30] H. Lim, D.-W. Kim, Pairwise dependence-based unsupervised feature  
552 selection, *Pattern Recognition* 111 (2021) 107663.
- 553 [31] P. Zhou, L. Du, X. Li, Y.-D. Shen, Y. Qian, Unsupervised feature se-  
554 lection with adaptive multiple graph learning, *Pattern Recognition* 105  
555 (2020) 107375.
- 556 [32] E. Jang, S. Gu, B. Poole, Categorical reparameterization with gumbel-  
557 softmax, in: *ICLR (Poster)*, OpenReview.net, 2017.
- 558 [33] K. Chaudhari, A. Thakkar, Survey on handwriting-based personality  
559 trait identification, *Expert Systems with Applications* 124 (2019) 282 –  
560 308.
- 561 [34] D. Castro, J. Tan, B. Kainz, E. Konukoglu, B. Glocker, Morpho-mnist:  
562 Quantitative assessment and diagnostics for representation learning,  
563 *JMLR* 20.
- 564 [35] P. Tschandl, C. Rosendahl, H. Kittler, The ham10000 dataset: A large  
565 collection of multi-source dermatoscopic images of common pigmented  
566 skin lesions, *Scientific Data* 5.
- 567 [36] D. Hovy, A. Johannsen, A. Sjøgaard, User review sites as a resource  
568 for large-scale sociolinguistic studies, in: *Proc. of WWW '15*, 2015, pp.  
569 452–461.