

Continual Rare-Class Recognition with Emerging Novel Subclasses

Hung Nguyen Xuejian Wang Leman Akoglu

Carnegie Mellon University

Heinz College of Information Systems and Public Policy

{hungnguy, xuejianw, lakoglu}@andrew.cmu.edu

Abstract. Given a labeled dataset that contains a rare (or minority) class of *of-interest* instances, as well as a large class of instances that are *not* of interest, how can we learn to recognize future *of-interest* instances over a continuous stream? We introduce RARECOGNIZE, which (i) estimates a *general* decision boundary between the rare and the majority class, (ii) learns to recognize individual rare subclasses that exist within the training data, as well as (iii) flags instances from previously unseen rare subclasses as newly emerging. The learner in (i) is general in the sense that by construction it is dissimilar to the *specialized* learners in (ii), thus distinguishes minority from the majority without overly tuning to what is seen in the training data. Thanks to this generality, RARECOGNIZE ignores all future instances that it labels as majority and recognizes the recurrent as well as emerging *rare* subclasses only. This saves effort at test time as well as ensures that the model size grows moderately over time as it only maintains specialized minority learners. Through extensive experiments, we show that RARECOGNIZE outperforms state-of-the art baselines on three real-world datasets that contain corporate-risk and disaster documents as rare classes.

1 Introduction

Given a labeled dataset containing (1) a rare (or minority) class of *of-interest* documents, and (2) a large set of *not-of-interest* documents, how can we learn a model that can effectively identify future *of-interest* documents over a continuous stream? Different from the traditional classification setup, the stream might contain *of-interest* (as well as *not-of-interest*) documents from *novel subclasses that were not seen in the training data*. Therefore, the model is required to continually recognize both the recurring as well as the emerging instances from the underlying rare class distribution.

Let us motivate this setting with a couple of real-world examples. Suppose we are given a large collection of social media documents (e.g. Twitter posts). A subset of the collection is labeled as *risky*, indicating posts that constitute (financial, reputational, etc.) risk to a corporation. The rest (majority) of the collection is *not-risky*. The goal is then to learn a model that can continually identify future posts that are *risky* over the social-media stream. Here, the rare class

contains *risky* documents of a few known types, such as bankruptcy, corruption, and spying. However, it is unrealistic to assume that it contains examples from all possible risk types—given the large spectrum, labeling effort, and potentially evolving nature of risk.

Consider another case where the training set consists of news articles. A subset of the articles belongs to the rare class of *disasters*, indicating news about natural or man-made disasters. The rest are *not-disaster* articles. Similar to the first case, the rare class might contain articles about floods, earthquakes, etc. however it is hard to imagine it would contain instances from all possible types of disasters. The goal is to learn to continually recognize future articles on disasters.

In both examples above, the model needs to learn from and generalize beyond the labeled data so as to recognize future rare-class instances, both from *recurring* (i.e., seen in the training data) as well as from *novel subclasses*; for instance sexual assault, cyber attack, etc. in risk domain and explosions, landslides, etc. in disasters domain. In machine learning terms, this is a very challenging setup in which the learner needs to generalize not only to unseen instances but also to *unseen distributions*. In other words, this setting involves test data that has a related yet different distribution than the data the model was trained on.

The stream classification problem under emerging novel classes has been studied by both machine learning and data mining communities. The area is referred to under various names including open-world classification [12,13], life-long learning [1], and continual learning [11]. In principle, these build a “never-ending learner” that can (1) assign those recurring instances from known old classes to their respective class, (2) recognize emerging classes, and (3) grow/extend the current learner to incorporate the new class(es). The existing methods differ in terms of accuracy-efficiency trade-offs and various assumptions that they make. (See Section 5 for detailed related work.) A common challenge that all of them face is what is known as *catastrophic forgetting*, mainly due to model growth. In a nutshell, the issue is the challenge of maintaining performance on old classes as the model is constantly grown to accommodate the new ones.

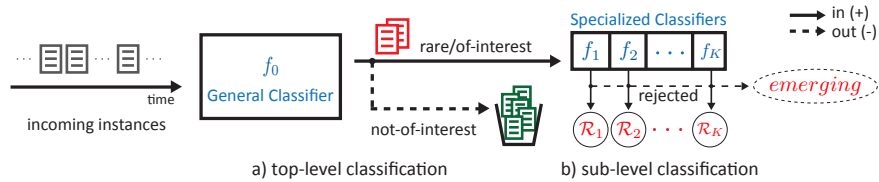


Fig. 1: An illustration of the recognition flow in our proposed model.

Our work is different from all prior work in one key aspect: our goal is not to recognize *any and every* newly emerging class—but only those (sub)classes related to the rare class *of-interest*. That is, our primary goal is to recognize rare-class instances. *Not-of-interest* instances, as long as they are filtered out accurately, are ignored—*no matter they are recurrent or novel*, as depicted in Fig. 1. This way, we carefully avoid the aforementioned issue that current models face. Our model grows slowly, only when novel *rare* subclasses are recognized.

Thanks to a moderate model size (by definition, rare subclasses are far fewer), our model is not only less prone to catastrophic forgetting but also (a) is faster at test time, and (b) requires much less memory.

We summarize the main contributions of this work as follows.

- **Problem and Formulation:** We address the problem of recognizing instances from a rare, *of-interest* class over a stream continually. The setting differs from traditional (binary) classification in that the data distribution (for both rare and majority class) might change over time, where novel subclasses emerge. We formulate a new model called RARECOGNIZE that *simultaneously* learns (i) a separate specialized classifier (SC) that recognizes an individual rare subclass, as well as (ii) a general classifier (GC) that separates rare instances from the majority. While being discriminative, GC is constructed to be dissimilar to the individual SCs such that it can generalize without overly tuning to seen rare subclasses in the training data.
- **Efficient Algorithm:** Our proposed solution exhibits two key properties: runtime and memory efficiency; both essential for the stream setting. Given a new instance that GC labels as belonging to the majority class, we simply do nothing—no matter it is recurrent or emerging. By not processing the majority of the incoming instances, we achieve *fast response time*. Moreover RARECOGNIZE remains compact, i.e. *memory-efficient*, as it requires space linear in the number of *rare* subclasses which only grows slowly.
- **Applications:** Recognizing recurrent as well as novel instances that belong to a certain class *of-interest* is a broad problem that finds numerous applications, e.g. in monitoring and surveillance. For example, such instances could be production-line items with the goal to continually recognize faulty ones where novel fault types might emerge over time. They could also be public documents, such as social media posts, where the goal is to recognize public posts *of-interest* such as bullying, shaming, disasters, threat, etc.

Reproducibility: We share the source code for RARECOGNIZE and our public-domain datasets at <https://github.com/hungnt55/RaRecognize>.

2 Problem Setup and Preliminary Data Analysis

Problem Setup and Overview. We start by introducing the problem statement more formally with proper notation. As input, a labeled training dataset $\mathcal{D} = \mathcal{R} \cup \mathcal{N} \in \mathbb{R}^{n \times d}$ containing n d -dimensional instances is provided. The set $\mathcal{R} = \{(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_{n_0}, y_{n_0})\}$ consists of $|\mathcal{R}| = n_0$ instances belonging to the *of-interest* rare class where $y_i = +1$ for $i = 1, \dots, n_0$ and the set $\mathcal{N} = \{(\mathbf{x}_{n_0+1}, y_{n_0+1}), \dots, (\mathbf{x}_n, y_n)\}$ consists of $|\mathcal{N}| = (n - n_0)$ instances from the *not-of-interest* class where $y_i = -1$ for $i = (n_0 + 1), \dots, n$. Without loss of generality, we will refer to the data instances as documents and to the rare class as the *risk* class in the rest of this section to present our ideas more concretely.

Given \mathcal{D} , the goal is to recognize future *risk* documents, *either recurring or newly emerging*, over a stream (or set) of new documents $\mathbf{x}_{n+1}, \mathbf{x}_{n+2}, \dots$ (here,

each document has a vector representation denoted by \mathbf{x} such as bag-of-words, embedding, etc.). The new documents may associate with recurring risk, i.e., belong to known/seen risk subclasses $1, \dots, K$ in \mathcal{R} . They may also be emerging, i.e., from previously unknown/unseen new risk subclasses $(K + 1), (K + 2), \dots$; which differentiates our setup from the traditional classification problem.

Therefore, we start by decomposing \mathcal{R} into known risk subclasses, $\mathcal{R} = \bigcup_{k=1}^K \mathcal{R}_k$, where \mathcal{R}_k contains the documents that belong to the k th risk subclass. Given $\{\{\mathcal{R}_1, \dots, \mathcal{R}_K\}, \mathcal{N}\}$ our approach involves simultaneously training the following two types of classifiers:

1. A *general* classifier (GC) f_0 to separate \mathcal{R} and \mathcal{N} that can generalize to unseen subclasses of \mathcal{R} ,
2. A *specialized* classifier (SC) f_k , $k = 1 \dots K$, to separate \mathcal{R}_k and $\mathcal{R} \setminus \mathcal{R}_k$.

At test time, we first employ f_0 . Our goal is not to recognize every emerging novel class, but only the novel risk subclasses (in addition to recurring ones), thus our first step is to recognize risk. If f_0 labels an incoming document \mathbf{x} as -1 (i.e., not-risk), we discard it. Otherwise, the incoming document is flagged as risky. For only those labeled as $+1$, we employ f_k 's to further identify the type of risk. Among the f_k 's that accept \mathbf{x} as belonging to the k th risk subclass, we assign it to the subclass that is $\arg \max_k f_k(\mathbf{x})$. If all f_k 's reject, then \mathbf{x} is considered to be associated with a new type of emerging risk. (See Fig. 1.)

The classifier models. Our risk detector is f_0 which we learn using the entire labeled dataset \mathcal{D} . As such, it is trained on a few known risk subclasses in \mathcal{R} but is desired to be *general* enough to recognize other types of future risk.

To achieve this generality, our main idea is to avoid building f_0 on factors that are too specific to any known risk subclass (such that f_0 is not overly fit to existing or known risk types) but rather, to identify broad factors about risk that are *common* to all risk subclasses (such that f_0 can employ this broader view to spot risk at large).

In fact, factors specific to the known risk subclasses are to be captured by the corresponding f_k 's. Then, f_0 is to identify discriminative signals of risk that are sufficiently different from those used by all f_k 's. Moreover, each f_k should differ from other $f_{k'}$'s, $k' \neq k$, to ensure that they are as *specialized* as possible to their respective risk types. Such dependence among the models is exactly why we train all these $(K + 1)$ classifiers *simultaneously*, to enforce the aforementioned constraints conjointly. We present our specific model formulation and optimization in Section 3.

Preliminary Data Analysis. Before model formulation, we perform an exploratory analysis on one of our real-world datasets containing documents labeled as risky and not-risky. The goal of the analysis is to see if our hypothesized ideas get realized in the data.

In particular, we aim to find out if there exists (1) factors that are specific to each risk subclass, as well as (2) factors beyond those specific ones that are still discriminative of risk. For simplicity and interpretability, we use the bag-of-words representation of the data in this section, thus factors correspond

to individual words. However, our proposed model can handle other document vector representations in general.

To this end, we formulate a constrained optimization problem to find word sets that cover or characterize different document sets. Here, we define a word to *cover* a document if the word appears in it at least once. Given the set of unique words \mathcal{V} , $|\mathcal{V}| = d$, we look for a set of words $\mathcal{V}_k \subset \mathcal{V}$ that covers all the documents in \mathcal{R}_k but as few as those in $\mathcal{R} \setminus \mathcal{R}_k$ for all $k = 1 \dots K$ (i.e., specific words for each risk subclass), and another set of words $\mathcal{V}_0 \subset \mathcal{V}$ that covers all risk documents in \mathcal{R} but as few as those in \mathcal{N} . We restrict the word sets to be *non-overlapping*, i.e., $\mathcal{V}_k \cap \mathcal{V}_{k'} = \emptyset \ \forall k, k' \in \{0, \dots, K\}$, such that each word can only characterize either one of $\mathcal{R}_1, \dots, \mathcal{R}_K$ or \mathcal{R} at large. Under these conditions, if we could find a set \mathcal{V}_0 that shares no words with any \mathcal{V}_k 's while still being able to cover the risky documents but only a few (if at all) not-risky ones, then we can conclude that broad risk terms exist and a *general* f_0 can be trained.

Our setup is a constrained mixed-integer linear program (MILP) as follows:

$$\begin{aligned}
\min_{\Phi} \quad & |\mathbf{v}_0| + \sum_{k=1}^K |\mathbf{v}_k| + o + \alpha + \beta \\
\text{s.t.} \quad & \mathbf{z}_k + \mathbf{R}_k \cdot \mathbf{v}_k \geq \mathbf{1} \quad \forall k = 1 \dots K &> \text{risk subclass coverage with exoneration} \\
& \mathbf{z}_0 + \sum_{k=1}^K \mathbf{R}_k \cdot \mathbf{v}_k \geq \mathbf{1} &> \text{risk coverage with exoneration} \\
& |\mathbf{z}_0| + \sum_{k=1}^K |\mathbf{z}_k| \leq o &> \# \text{ unexplained documents less than } o \\
& \mathbf{v}_0 + \sum_{k=1}^K \mathbf{v}_k \leq \mathbf{1} &> \text{each word used for at most 1 set} \\
& \sum_{i \in \mathcal{N}} \mathbf{N}_i \cdot \mathbf{v}_0 \leq \alpha &> \text{cap on cross-coverage of not-risk} \\
& \sum_{k=1}^K \sum_{i \in \mathbf{R}_k} \sum_{k' \neq k} \mathbf{R}_{k,i} \cdot \mathbf{v}_{k'} \leq \beta &> \text{cap on cross-coverage among risk subclasses}
\end{aligned}$$

The program is parameterized by $\Phi = \{\{\mathbf{v}_k\}_{k=0}^K, \{\mathbf{z}_k\}_{k=0}^K, o, \alpha, \beta\}$. $\mathbf{R}_k \in \mathbb{R}^{n_k \times d}$ denotes the data matrix encoding the word occurrences for n_k documents in risk subclass $k = 1 \dots K$, and $\mathbf{N} \in \mathbb{R}^{(n-n_0) \times d}$ is the corresponding data matrix for the not-risk documents. $\mathbf{v}_0 \in \mathbb{R}^d$ and $\mathbf{v}_k \in \mathbb{R}^d$'s depict (binary) variables to be estimated that capture the word assignments to the sets \mathcal{V}_0 and \mathcal{V}_k 's respectively. (e.g., j th entry of \mathbf{v}_0 is set to 1 if word j is assigned to \mathcal{V}_0 and to 0 otherwise.)

The first set of constraints are coverage constraints for risk subclasses: each document in \mathcal{R}_k should contain at least one of the words in their assigned set. Enforcing this constraint for *all* the documents is too strict, hence we introduce additional (binary) variables $\mathbf{z}_k \in \mathbb{R}^{n_k}$'s that ‘‘exonerate’’ some documents. When i th entry of \mathbf{z}_k is set to 1, then document i in \mathcal{R}_k is allowed to have no matching words, as $\mathbf{z}_k = \mathbf{1}$ ensures the constraint holds even without any match. The second constraint is similar, and enforces coverage for the combined set of risk

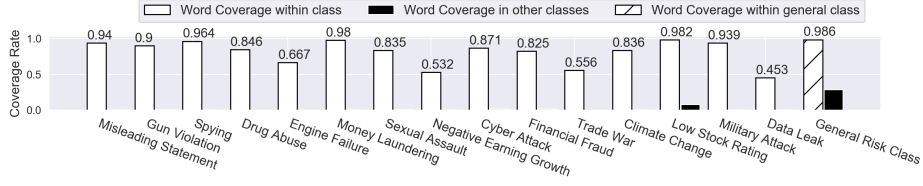


Fig. 2: Within- and cross-coverage rates of \mathcal{V}_k 's and \mathcal{V}_0 (resp.) for \mathcal{R}_k 's and \mathcal{R} .

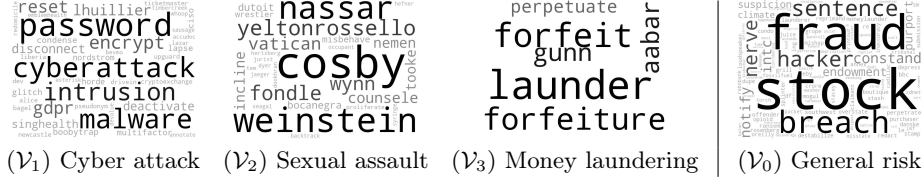


Fig. 3: Wordclouds representing 3 example risk subclasses and the overall risk class. Notice that the former are quite specific, and the latter are broader.

documents, also with exoneration. Of course, we aim to cover as many documents as possible and thus upper-bound the total number of exonerated documents by o , where $|\mathbf{z}|$ denotes the total number of 1s in vector \mathbf{z} . Next is the no-overlaps constraint, enforcing each word is assigned to only one set. The final two constraints are cross-coverage constraints; the former ensures that the words assigned to \mathcal{V}_0 have less than α number of matches in not-risk documents and the latter ensures that the words assigned to each risk subclass have less than β number of matches outside the respective document set in total.

Ideally all of o , α , and β are zero; that is, all documents are covered without any exoneration and no cross-coverage exists. However, that yields no feasible solution. Instead, we define them as scalar upper-bound variables added to our minimization objective toward setting them to as small values as possible. Finally, our objective aims to find the smallest-size possible word sets. This ensures that the most important words are selected which also facilitates interpretability.

We provide an exploratory analysis on a dataset containing corporate risk documents as the *of-interest* class. It contains 15 risk subclasses as outlined by Fig. 2. (See Sec. 4.1 for details.) First the quantitative measures: as shown in the figure, the MILP finds word sets \mathcal{V}_k 's with at least 82.5% up to 98.2% coverage for 11/15 of the subclasses with an overall coverage of 96.7% (rest are the exonerated ones). Moreover, cross-coverage is either zero or very low for all the subclasses. These suggest that accurate SCs can be learned. Importantly, there exist words \mathcal{V}_0 that are *distinct* from all \mathcal{V}_k 's and yet able to cover 98.6% of the overall risk documents, promising that a broad GC can be learned.

To equip the reader with intuition, we present the selected words for 3 example risk subclasses along with the general risk class words in Fig. 3 (word size is proportional to the within vs. cross-coverage ratio). It is easy to see that

very specific words are selected for subclasses; such as **password**, **cyberattack**, **malware** for Cyber attack, and **cosby**, **weinstein**, **fondle** for Sexual assault. On the other hand, in the general risk class, a set of broader corporate risk words appear such as **fraud**, **stock**, **breach** and **sentence**.

These preliminary results show promise for the feasibility of our hypothesized models and demonstrate the rationale behind our proposed RARECOGNIZE, which we formally introduce next.

3 Continual Rare-Class Recognition

In this section, we introduce the individual components of our model, present the underlying reasoning for our formulation, show convexity and present the optimization steps, and conclude with space and time-complexity analysis.

3.1 Model Formulation

As discussed in the previous section, our goal is to learn (1) specialized classifiers f_k 's and (2) a general classifier f_0 .

The **specialized classifier** f_k , $k = 1 \dots K$, is to learn a decision boundary that separates the k th rare subclass instances \mathcal{R}_k from the remainder of rare instances $\mathcal{R} \setminus \mathcal{R}_k$. Let us write down the regularized loss function for each f_k as

$$\mathcal{L}(f_k; \mathbf{w}_k, b_k) = \sum_{i=1}^{n_0} \underbrace{\max(0, [1 - y_i(\mathbf{w}_k^T \mathbf{x}_i + b_k)])}_{\ell(\mathbf{x}_i, y_i; \mathbf{w}_k, b_k)} + \frac{\lambda_k}{2} \|\mathbf{w}_k\|^2 \quad (1)$$

where $y_i = +1$ for $\mathbf{x}_i \in \mathcal{R}_k$ and $y_i = -1$ otherwise. We adopt the hinge loss and the ridge regularization as in Eq. (1), however, one could instead use other loss functions, such as the logistic, exponential or cross-entropy losses, as well as other norms for regularization.

The **general classifier** f_0 is to separate rare class instances \mathcal{R} from the majority instances \mathcal{N} , without relying on factors specific to known rare subclasses. One way to achieve this *de-correlation* is to enforce f_0 to learn coefficients \mathbf{w}_0 that are different from all \mathbf{w}_k 's. The loss function can be written as

$$\mathcal{L}(f_0; \mathbf{w}_0, b_0) = \sum_{i=1}^n \ell(\mathbf{x}_i, y_i; \mathbf{w}_0, b_0) + \frac{\lambda_0}{2} \|\mathbf{w}_0\|^2 + \frac{\mu}{2} \sum_{k=1}^K \mathbf{w}_0^T \mathbf{w}_k, \quad (2)$$

where $y_i = +1$ for $\mathbf{x}_i \in \mathcal{R} = \{\mathcal{R}_1 \cup \dots \cup \mathcal{R}_K\}$ and $y_i = -1$ otherwise. As required, the third term in Eq. (2) penalizes \mathbf{w}_0 being correlated with any \mathbf{w}_k , enforcing it to be as orthogonal to \mathbf{w}_k 's as possible. However, it does not prevent \mathbf{w}_0 from capturing *different yet correlated features* to those captured by \mathbf{w}_k 's. This issue can arise when features exhibit multi-collinearity.

For example, in a document dataset the words **earthquake**, **shockwave**, and **aftershock** could be collinear. In this case it is possible that f_k estimates large coefficients on a strict subset of these words (e.g., **shockwave** and **aftershock**) as they are redundant. This leaves room for f_0 to capitalize on the remaining

words (e.g., **earthquake**), which is undesirable since we aim f_0 to learn about the rare class boundaries beyond the specifics of the known subclasses.

Therefore, we reformulate the model correlation penalty as follows.

$$\frac{\mu}{2} \sum_{k=1}^K \sum_{p,q} (w_{0,p} w_{k,q} \mathbf{x}_{[p]}^T \mathbf{x}_{[q]})^2 = \frac{\mu}{2} \|(\mathbf{X}^T \mathbf{X}) \odot (\mathbf{w}_0 \mathbf{w}_k^T)\|_F^2, \quad (3)$$

where $w_{0,p}$ and $w_{k,q}$ denote the p th and q th entries of \mathbf{w}_0 and \mathbf{w}_k respectively, $\mathbf{X} \in \mathbb{R}^{n \times d}$ denotes the input data matrix, $\mathbf{x}_{[p]}, \mathbf{x}_{[q]}$ respectively denote the p th and q th columns of \mathbf{X} , and \odot depicts the element-wise multiplication.

We call Eq. (3) the cross-correlation penalty. Similarly, we introduce self-correlation penalty to each model $k = 0, \dots, K$ by adding to the respective loss the term $\sum_{p,q} (w_{k,p} w_{k,q} \mathbf{x}_{[p]}^T \mathbf{x}_{[q]})^2$. Self-correlation prevents each model from estimating large coefficients on highly correlated (near-redundant) features, which improves sparsity and interpretability, and as we show also ensures convexity.

Then, the overall loss function incorporating the cross- and self-correlation penalty terms for all models f_0, f_1, \dots, f_K is given as follows.

$$\mathcal{L} = \sum_{i=1}^n \ell(\mathbf{x}_i, y_i; \mathbf{w}_0, b_0) + \frac{\lambda_0}{2} \|\mathbf{w}_0\|^2 + \sum_{k=1}^K \left[\sum_{i=1}^{n_0} \ell(\mathbf{x}_i, y_i; \mathbf{w}_k, b_k) + \frac{\lambda_k}{2} \|\mathbf{w}_k\|^2 \right] + \frac{\mu}{2} \sum_{p,q} \left\{ \underbrace{\frac{1}{2} (w_{0,p}^2 w_{0,q}^2)}_{\text{self-correlation}} + \underbrace{\frac{1}{2} \sum_{k=1}^K (w_{k,p}^2 w_{k,q}^2) + w_{0,p}^2 \left(\sum_{k=1}^K w_{k,q}^2 \right)}_{\text{cross-correlation}} \right\} (\mathbf{x}_{[p]}^T \mathbf{x}_{[q]})^2 \quad (4)$$

3.2 Convexity and Optimization

We train all the models f_0, f_1, \dots, f_K *simultaneously* by minimizing the total overall loss \mathcal{L} . A conjoint optimization is performed because the cross-correlation penalty terms between \mathbf{w}_0 and \mathbf{w}_k 's induce dependence between the models.

For optimization we employ the accelerated subgradient descent algorithm which is guaranteed to find the global optimum solution because, as we show next, our loss function \mathcal{L} is convex.

Theorem 1. *The joint loss function \mathcal{L} involving the cross- and self-correlation penalty terms among $\mathbf{w}_0, \mathbf{w}_1, \dots, \mathbf{w}_K$ remains convex.*

Proof. The non-negative sum of convex functions is also convex. The first line of \mathcal{L} as given above is known to be convex since $\ell(\cdot)$ (hinge loss) and L- p norms for $p \geq 1$ are both convex. The proof is then by showing that the overall correlation penalty term in the second line of \mathcal{L} is also convex by showing that its Hessian matrix is positive semi-definite (PSD). See Supplementary A.1 for details. \square

Since our total loss is a convex function, we can use gradient-based optimization to solve it to optimality. To this end, we provide the gradient updates for both \mathbf{w}_0 and \mathbf{w}_k 's in closed form as follows.

Partial derivative of \mathcal{L} w.r.t. \mathbf{w}_0 :

$$\frac{\partial \mathcal{L}}{\partial w_{0,p}} = \sum_{i=1}^n \frac{\partial [1 - y_i(\mathbf{w}_0^T \mathbf{x}_i + b_0)]_+}{\partial w_{0,p}} + \lambda_0 w_{0,p} + \mu w_{0,p} \sum_{q=1}^d \left(\sum_{k=1}^K w_{k,q}^2 + w_{0,q}^2 \right) (\mathbf{x}_{[p]}^T \mathbf{x}_{[q]})^2$$

where

$$\frac{\partial [1 - y_i(\mathbf{w}_0^T \mathbf{x}_i + b_0)]_+}{\partial w_{0,p}} = \begin{cases} 0 & \text{if } y_i(\mathbf{w}_0^T \mathbf{x}_i + b_0) \geq 1 \\ -y_i x_{i,p} & \text{otherwise.} \end{cases} \quad (5)$$

The vector-update $\frac{\partial \mathcal{L}}{\partial \mathbf{w}_0}$ can be given in matrix-vector form using the above gradients as

$$\frac{\partial \mathcal{L}}{\partial \mathbf{w}_0} = [\mathbf{X}^T (-\mathbf{y} \odot \mathbb{I}(1 - \mathbf{y} \odot (\mathbf{X} \mathbf{w}_0 + b_0) > 0)) + \mathbf{w}_0 \odot (\lambda_0 \mathbf{1} + \mu (\mathbf{X}^T \mathbf{X})^2 (\sum_k \mathbf{w}_k^2 + \mathbf{w}_0^2))] \quad (6)$$

where $\mathbb{I}(\cdot)$ is the indicator function, $\mathbf{1}$ is a length- n all-ones vector, and $(\mathbf{A})^2 = \mathbf{A} \odot \mathbf{A}$, i.e., element-wise product, for both matrix \mathbf{A} as well as for vector \mathbf{a} .

Partial derivative of \mathcal{L} w.r.t. \mathbf{w}_k : The steps for each \mathbf{w}_k is similar, we directly provide the vector-update below.

$$\frac{\partial \mathcal{L}}{\partial \mathbf{w}_k} = [\mathbf{R}^T (-\mathbf{y}_0 \odot \mathbb{I}(1 - \mathbf{y}_0 \odot (\mathbf{R} \mathbf{w}_k + b_k) > 0)) + \mathbf{w}_k \odot (\lambda_k \mathbf{1} + \mu (\mathbf{X}^T \mathbf{X})^2 (\mathbf{w}_k^2 + \mathbf{w}_0^2))] \quad (7)$$

where $\mathbf{R} \in \mathbb{R}^{n_0 \times d}$ and $\mathbf{y}_0 \in \mathbb{R}^{n_0}$ consist of only the rare-class instances.

3.3 Time and Space-Complexity Analysis

Time Complexity. The (first) gradient term in Eq (6) that is related to the hinge-loss is $\mathcal{O}(nd)$. The (second) term related to the correlation-based regularization requires $\mathbf{X}^T \mathbf{X}$ which can be computed in $\mathcal{O}(nd^2)$ apriori and *reused* over the gradient iterations. The term $(\sum_k \mathbf{w}_k^2 + \mathbf{w}_0^2)$ takes $\mathcal{O}(Kd)$, and its following multiplication with $(\mathbf{X}^T \mathbf{X})^2$ takes an additional $\mathcal{O}(d^2)$. The remaining operations (summation with $\lambda_0 \mathbf{1}$ and element-wise product with \mathbf{w}_0) are only $\mathcal{O}(d)$. As such, the overall computational complexity for subgradient descent for \mathbf{w}_0 is $\mathcal{O}(nd^2 + t[nd + Kd + d^2])$, where t is the number of gradient iterations.

The time complexity for updating the \mathbf{w}_k 's can be derived similarly as Eq. (7) consists of similar terms, which can be written as $\mathcal{O}(t[n_0 d + d^2])$. Note that we omit the $\mathcal{O}(nd^2)$ this time as $\mathbf{X}^T \mathbf{X}$ needs to be computed only once and can be shared across all update rules. Moreover, the number of iterations t is the same as before since the parameter estimation is conjoint.

Space Complexity. We require $\mathcal{O}(d^2)$ storage for keeping $(\mathbf{X}^T \mathbf{X})^2$, $\mathcal{O}(Kd)$ for all the parameter vectors $\mathbf{w}_0, \dots, \mathbf{w}_K$, and $\mathcal{O}(nd)$ for storing \mathbf{X} , for an overall $\mathcal{O}(d^2 + Kd + nd)$ space complexity.

Remarks on massive and/or high-dimensional datasets: Note that both time and space complexity of our RARECOGNIZE is quadratic in d and linear in n . We conclude with parting remarks on cases with large d and huge n .

First, high-dimensional data with large d : In this case, we propose two possible directions to make the problem tractable. Of course, the first one is dimensionality reduction or representation learning. When the data lies on a relatively

low-d manifold, one could instead use compound features. We apply our RARECOGNIZE to document datasets, where compound features are not only fewer but also sufficiently expressive of the data. The second direction is to get rid of feature correlations, for instance via factor analysis. This would drop the term $(\mathbf{x}_{[p]}^T \mathbf{x}_{[q]})^2$ from \mathcal{L} (See (4)) and lead to updates that are only *linear* in d .

Next, massive data with huge n : We presented our optimization using batch subgradient descent. When n is very, very large then storing the original data in memory may not be feasible. We remark that one could directly employ mini-batch or even stochastic gradient descent in such cases, dropping the space requirement to $O(d^2 + Kd)$.

4 Evaluation

We design experiments to evaluate our proposed method with respect to the following questions:

- **RQ1) Top-level classification (via GC f_0):** How does RARECOGNIZE perform in differentiating rare-class instances from the majority compared with the state-of-the-art?
- **RQ2) Sub-level classification (via SCs f_k ’s):** How does RARECOGNIZE perform in recognizing recurrent and emerging rare subclasses among the compared methods?
- **RQ3) Interpretability:** Can we interpret RARECOGNIZE as a model as to what it has learned and what insights can we draw?
- **RQ4) Efficiency:** What is the scalability of RARECOGNIZE? How does it compare to the baselines w.r.t. the running time-vs-performance trade-off?

4.1 Experiment Setup

Dataset Description. In this study, we use 3 different datasets with characteristics summarized in Table 1. The first two datasets are obtained from our industry collaborator (proprietary) and a third public one which we put together.

RISK-DOC: This dataset contains online documents, e.g. news articles, social network posts, which are labeled risky or non-risky to the corporate entities mentioned. If a document is risky, it is further assigned to one of the 15 risk subclasses: {Climate change, Cyber attack, Data leak, Drug abuse, Engine failure, Fraud, Gun violation, Low stock, Military attack, Misleading statement, Money laundering, Negative growth, Sexual assault, Spying, Trade war}.

RISK-SEN: This contains labeled sentences attracted from news articles and categorized into 8 different subclasses: {Bankruptcy, Bribery corruption, Counterfeiting, Cyber privacy, Environment, Fraud false claims, Labor, Money laundering}. The majority class consists of non-risky sentences. Note that this dataset comes from a set of articles different from RISK-DOC.

Table 1: Summary of datasets.

Name	$ \mathcal{R} $	$ \mathcal{N} $	K	Avg. \mathcal{R}_k
RISK-DOC	2948	2777	15	196.5
RISK-SEN	1551	7755	8	193.9
NYT-DSTR	2127	10560	13	163.6

NYT-DSTR: Extracted from the New York Times, this dataset is comprised of articles on the topics of disasters, i.e. both natural and human-instigated. These topics cover 13 disasters from {Drought, Earthquakes, Explosions, Floods, Forest and bush fire, Hazardous and toxic substance, Landslides, Lighting, Snowstorms, Tornado, Tropical storm, Volcanoes, Water pollution}. It also includes a class of random non-disaster news articles from New York Times.

Document representations. In this study we apply our work to document datasets, for which we need to define a feature representation. There are numerous options. We report results with tfidf with top 1K words based on frequency, as well as PCA- and ICA-projected data. Linear embedding techniques reduce dimensionality while preserving interpretability. We omit results using non-linear feature representations (e.g., doc2vec [6]) as they did not provide any significant performance gain despite computational overhead.

Train/Test Splits. For each dataset, we randomly partition 2/3 of rare subclasses as seen and 1/3 of them as unseen. For training, we use 80% of the seen subclass instances at random and the rest 20% forms a seen subclass test set, denoted by \mathcal{R}_s . The set of unseen subclass instances, denoted by \mathcal{R}_u , goes into the test as well. Thus, the rare subclass test consists of 2 parts, i.e. $\mathcal{R}_{\text{test}} = \mathcal{R}_s \cup \mathcal{R}_u$. In addition, we also reserve a random 80% of the majority class for training and the rest 20% for testing, denoted by $\mathcal{N}_{\text{test}}$. Further, to obtain stable results, we repeat our experiment on 5 different random train/test constructions and report averaged outcomes.

Performance Metrics. (1) For measuring top-level classification performance, we use 3 common metrics [15,8]: *Precision*, *Recall*, *F1* formally defined in our context as:

$$Precision = \frac{|\mathcal{R}_{\text{test}} \cap \hat{\mathcal{R}}_{\text{test}}|}{|\hat{\mathcal{R}}_{\text{test}}|}, Recall = \frac{|\mathcal{R}_{\text{test}} \cap \hat{\mathcal{R}}_{\text{test}}|}{|\mathcal{R}_{\text{test}}|}, F1 = \frac{2 * Precision * Recall}{Precision + Recall},$$

where $\hat{\mathcal{R}}_{\text{test}}$ is the set of examples predicted as rare subclass. To identify which part of the test (seen or unseen subclasses) the model makes mistakes, we also measure *Precision (seen)*, *Recall (seen)* and *Recall (unseen)* defined as follows:

$$Precision(\text{seen}) = \frac{|\mathcal{R}_s \cap \hat{\mathcal{R}}_{\text{test}}|}{|\hat{\mathcal{R}}_{\text{test}}|}, Recall(\text{seen}) = \frac{|\mathcal{R}_s \cap \hat{\mathcal{R}}_{\text{test}}|}{|\mathcal{R}_s|}, Recall(\text{unseen}) = \frac{|\mathcal{R}_u \cap \hat{\mathcal{R}}_{\text{test}}|}{|\mathcal{R}_u|}.$$

(2) For sub-level classification test, to quantify the fraction of seen subclass test instances correctly classified and unseen subclass test instances as emerging, we use the following metric:

$$acc(\text{rare}) = \frac{\overbrace{|\mathcal{R}_u \cap \hat{\mathcal{R}}_u|}^{acc(\text{emerging})} + \sum_{k=1}^K \overbrace{|\mathcal{R}_{ks} \cap \hat{\mathcal{R}}_{ks}|}^{acc(\text{recurrent})}}{|\mathcal{R}_{\text{test}}|},$$

where \mathcal{R}_{ks} is the set of test examples in subclass k and $\hat{\mathcal{R}}_{ks}$ is the set of examples assigned to that subclass. Here $acc(\text{rare}) = 1$ if both seen subclass test instances are perfectly classified to their respective subclasses and unseen subclass instances as emerging. For all of the above metrics, the higher is better.

Compared Methods. We compare RARECOGNIZE with 2 state-of-the-art methods and 2 simple baselines:

- RARECOGNIZE-1K, RARECOGNIZE-PCA, RARECOGNIZE-ICA: 3 versions of RARECOGNIZE when tfidf with 1K word dictionary, PCA and ICA representations are used. In RARECOGNIZE-PCA, we drop the feature correlation terms $(\mathbf{x}_{[p]}^T \mathbf{x}_{[q]})^2$ since features are orthogonal. For the sub-level classification, RARECOGNIZE learns a rejection threshold for each specialized classifier based on extreme value theory [14].
- L2AC [15]: the most recent method (2019) in open-world classification setting that is based on deep neural networks. We use the recommended parameters $k = 5, n = 9$ (in their notation) from the paper.
- SENCFOREST [8]: another state-of-the-art ensemble method (2017) using random decision trees for classification under emerging classes. We run SENCFOREST with 100 trees and subsample size 100 as suggested in their paper.
- BASELINE: a baseline of RARECOGNIZE when both cross- and self-correlation terms in Eq. (4) are removed, via setting $\mu = 0$. Basically, BASELINE is independently trained f_0, \dots, f_K .
- BASELINE-r: a variant of BASELINE when classification threshold (0.5 by default) is chosen so that the *Recall* matches that of RARECOGNIZE.

Note that SENCFOREST and L2AC aim to detect *any* emerging class without categorizing into rare or majority. For fair comparison, we only include new *rare* subclasses in our test data and consider their rejected instances as belonging to those. In reality, however, emerging classes need to be categorized as rare or not, which these existing methods did not address.

Table 2: Performance of methods on the three datasets.

Methods	Precision			Recall			F1		
	RISK-DOC	RISK-SEN	NYT-DSTR	RISK-DOC	RISK-SEN	NYT-DSTR	RISK-DOC	RISK-SEN	NYT-DSTR
SENCFOREST	0.46±0.12	0.14±0.03	0.36±0.09	0.59±0.09	0.41±0.08	0.39±0.10	0.52±0.11	0.21±0.04	0.37±0.06
L2AC	0.79±0.08	0.47±0.06	0.31±0.07	0.57±0.29	0.85±0.05	0.76±0.07	0.63±0.24	0.60±0.04	0.44±0.07
BASELINE	0.89±0.04	0.79±0.05	0.86±0.03	0.52±0.06	0.55±0.04	0.63±0.03	0.65±0.05	0.65±0.03	0.73±0.03
BASELINE-r	0.76±0.10	0.56±0.07	0.71±0.07	0.58±0.13	0.57±0.14	0.79±0.04	0.65±0.10	0.55±0.06	0.75±0.04
RARECOGNIZE-1K	0.89±0.02	0.83±0.09	0.84±0.03	0.58±0.13	0.57±0.14	0.79±0.04	0.70±0.09	0.66±0.08	0.81±0.02
RARECOGNIZE-PCA	0.85±0.06	0.79±0.10	0.84±0.10	0.58±0.17	0.71±0.17	0.80±0.07	0.68±0.14	0.73±0.09	0.81±0.01
RARECOGNIZE-ICA	0.74±0.07	0.72±0.11	0.84±0.12	0.73±0.24	0.85±0.18	0.82±0.08	0.71±0.09	0.78±0.07	0.81±0.04

4.2 Experiment Results

In the following, we sequentially answer the questions by analyzing our experimental results and comparing between methods.

RQ1) Top-level Classification into Rare vs. Majority Class. We report the *Precision*, *Recall* and *F1* of all methods on three datasets in Table 2. For SENCFOREST, L2AC, BASELINE and BASELINE-r, we report results for the representation (tfidf top-1K, PCA, ICA) that yielded the highest *F1* value.

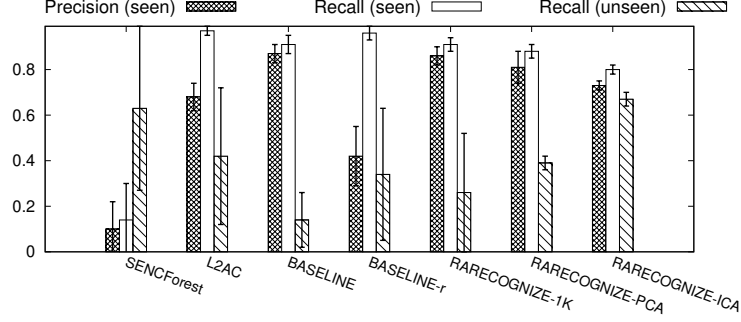


Fig. 4: *Precision (seen)*, *Recall (seen)* and *Recall (unseen)* of methods on RISK-DOC (RARECOGNIZE-ICA achieves the best balance between *Precision* and *Recall* on both seen and unseen test instances).

From Table 2, we see that RARECOGNIZE-1K, RARECOGNIZE-PCA and RARECOGNIZE-ICA outperform other methods in terms of *F1* score in all cases and *Precision*, *Recall* in most cases. Compared to BASELINE and BASELINE-r, *F1* score of RARECOGNIZE-ICA is 6-13% higher than the best result among the two. This demonstrates that cross- and self-correlations are crucial in RARECOGNIZE. Surprisingly, the gap to SENCFOREST and L2AC is even larger in terms of *F1*, between 8-37% higher. This shows that previous methods on detecting any new emerging classes do not work well when we only target rare subclasses.

Among the three versions of RARECOGNIZE, RARECOGNIZE-ICA gives the highest *F1*. RARECOGNIZE-ICA achieves the best balance between precision and recall while RARECOGNIZE-1K and RARECOGNIZE-PCA seem to have very high *Precision* but much lower *Recall*. That means that RARECOGNIZE-1K and RARECOGNIZE-PCA are better than RARECOGNIZE-ICA at discarding majority samples and worse at recognizing rare subclasses.

In Fig. 4, we have the *Precision (seen)*, *Recall (seen)* and *Recall (unseen)* measures of all the methods on RISK-DOC (Figures for other datasets are similar, see Supp. A.2). This figure shows that RARECOGNIZE-ICA also achieves a good balance between seen and unseen subclass classification, i.e., it recognizes both these subclasses equally well. On the other hand, most of other methods achieve high *Precision (seen)* and *Recall (seen)* but much lower *Recall (unseen)*, except SENCFOREST which only has high *Recall (unseen)*. This is because SENCFOREST rejects most instances as unseen which however hurts *Precision* drastically.

RQ2) Sub-level Classification into Recurrent and Emerging Rare Subclasses. We report the *acc(rare)* of all the methods in Table 3 (Breakdown of errors in confusion tables are given in Supp. A.3).

Tables 2 and 3 reflect that all three versions of RARECOGNIZE are always better or comparable to the others in terms of *acc(rare)* and RARECOGNIZE-ICA achieves the highest value. RARECOGNIZE-ICA achieves significantly higher *acc(rare)* than all the baselines. SENCFOREST seems to perform the next best

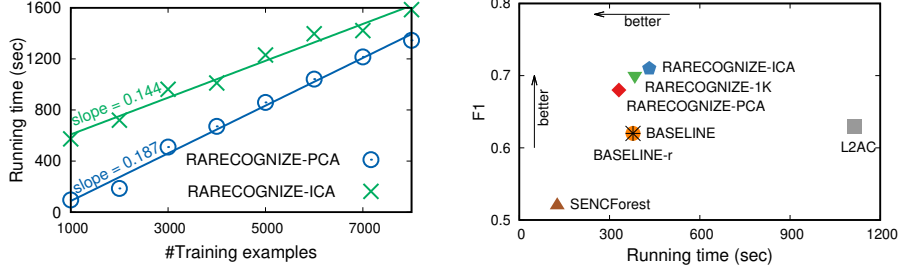


Fig. 6: RARECOGNIZE scales linearly in Fig. 7: Time-Performance trade-off of all number of training examples. compared methods.

early with the data size. RARECOGNIZE with PCA is faster than that with ICA thanks to no feature correlations (i.e. $(\mathbf{x}_{[p]}^T \mathbf{x}_{[q]})^2$ dropped).

In Fig. 7, we show the time-performance trade-off among all compared methods. We conclude that RARECOGNIZE with three representations run relatively fast, only slower than SENCFOREST, and returns the highest performance in terms of $F1$. L2AC consumes a huge amount of time for training a neural network, with subpar performance.

5 Related Work

Our work is closely related with two fields, namely open-world classification and continual learning. Both belong to the category of lifelong machine learning [1].

Open-world classification. Traditional close-world classification assumes that all test classes are known and seen in training data [4, 16]. However, such assumption could be violated in reality. Open-world classification, in contrast, assumes unseen and novel classes could emerge during test time, and addresses the classification problem by recognizing unseen classes. Previous works [12, 13, 15, 9, 8] propose different approaches under this setting.

Specifically, DOC [12] leverages convolutional neural nets (CNNs) with multiple sigmoid functions to classify examples as seen or emerging. [13] follows the same DOC module and performs hierarchical clustering to all rejected samples. Later, L2AC [15] proposes to use a meta-classifier and a ranker to add or delete a class without re-training. However, it requires a large amount of computation in both training and testing due to the top- k search over all training data.

SENCForest [8] is a randomized ensemble method. It grows multiple random forests and rejects examples when all random forests yield "new class". Under the same setting, SENC-MaS [9] maintains matrix sketchings to decide whether an example belongs to a seen class or emerging.

In the emerging rare subclass setting, the previous approaches aim at recognizing *any and every* classes and are not able to ignore the *not-of-interest* classes while recognizing emerging ones, thus consume much more memory and time.

Continual learning. There are recent works investigating continual learning or incremental learning [11,10]. They aim at solving the issue of catastrophic forgetting [2] in connectionist networks. In this field, models are proposed to continually learn new classes without losing performance on old seen classes.

Previous works [10,5,7,3] show promising results. However, the number of documents in rare subclasses *of-interest* in our setting is usually not large enough for neural networks to be sufficiently trained. Consequently, the neural network approach does not perform well in rare-class classification and recognition.

6 Conclusion

We proposed RARECOGNIZE for rare-class recognition over a continuous stream, in which new subclasses may emerge. RARECOGNIZE employs a general classifier to filter out not-rare class instances (top-level) and a set of specialized classifiers that recognize known rare subclasses or otherwise reject as emerging (sub-level). Since majority of incoming instances are filtered out and new rare subclasses are a few, RARECOGNIZE processes incoming data fast and grows in size slowly. Extensive experiments show that it outperforms two most recent state of the art as well as two simple baselines significantly in both top- and sub-level tasks, while achieving the best efficiency-performance balance and offering interpretability. Future work will extend RARECOGNIZE to an end-to-end system that clusters emerging instances and trains all the relevant models incrementally.

References

1. Z. Chen and B. Liu. Lifelong machine learning. *Synthesis Lectures on Artificial Intelligence and Machine Learning*, 10(3):1–145, 2016.
2. R. French. Catastrophic forgetting in connectionist networks. *Trends in cognitive sciences*, 3:128–135, 05 1999.
3. R. Kemker and C. Kanan. Fearnnet: Brain-inspired model for incremental learning. In *ICLR*, 2018.
4. Y. Kim. Convolutional neural networks for sentence classification. *EMNLP*, 2014.
5. J. Kirkpatrick, R. Pascanu, N. Rabinowitz, J. Veness, G. Desjardins, A. A. Rusu, K. Milan, J. Quan, T. Ramalho, A. Grabska-Barwinska, et al. Overcoming catastrophic forgetting in neural networks. *PNAS*, 114(13):3521–3526, 2017.
6. Q. V. Le and T. Mikolov. Distributed representations of sentences and documents. In *ICML*, volume 14, pages 1188–1196, 2014.
7. S.-W. Lee, J.-H. Kim, J. Jun, J.-W. Ha, and B.-T. Zhang. Overcoming catastrophic forgetting by incremental moment matching. In *NeurIPS*, pages 4652–4662, 2017.
8. X. Mu, K. M. Ting, and Z.-H. Zhou. Classification under streaming emerging new classes: A solution using completely-random trees. *IEEE TKDE*, 29(8), 2017.
9. X. Mu, F. Zhu, J. Du, E.-P. Lim, and Z.-H. Zhou. Streaming classification with emerging new class by class matrix sketching. In *AAAI*, 2017.
10. S.-A. Rebuffi, A. Kolesnikov, G. Sperl, and C. H. Lampert. icarl: Incremental classifier and representation learning. In *CVPR*, pages 2001–2010, 2017.
11. H. Shin, J. K. Lee, J. Kim, and J. Kim. Continual learning with deep generative replay. In *NeurIPS*, pages 2990–2999, 2017.

12. L. Shu, H. Xu, and B. Liu. Doc: Deep open classification of text documents. *EMNLP*, 2017.
13. L. Shu, H. Xu, and B. Liu. Unseen class discovery in open-world classification. *arXiv preprint arXiv:1801.05609*, 2018.
14. A. Siffer, P.-A. Fouque, A. Termier, and C. Largouet. Anomaly detection in streams with extreme value theory. In *KDD*, pages 1067–1075. ACM, 2017.
15. H. Xu, B. Liu, L. Shu, and P. Yu. Open-world learning and application to product classification. In *WWW*, 2019.
16. X. Zhang, J. Zhao, and Y. LeCun. Character-level convolutional networks for text classification. In *NeurIPS*, pages 649–657, 2015.

A Supplementary

A.1 Proof of Theorem 1.

Given that $\ell(\cdot)$ and L- p norms for $p \geq 1$ are convex, and that sum of non-negative convex functions remains convex, it suffices to show that the correlation term (denoted by C) of the loss function \mathcal{L} is convex.

$$C = \sum_{p,q} \left\{ \frac{\mu}{4} (w_{0,p}^2 w_{0,q}^2) + \frac{\mu}{4} \sum_{k'=1}^K (w_{k',p}^2 w_{k',q}^2) + \frac{\mu}{2} w_{0,p}^2 \left(\sum_{k'=1}^K w_{k',q}^2 \right) \right\} (\mathbf{x}_{[p]}^T \mathbf{x}_{[q]})^2 \quad (8)$$

We prove convexity by showing that the Hessian matrix of C is positive semi-definite (PSD). The Hessian in this case is a $d(K+1) \times d(K+1)$ matrix, denoted \mathbf{H} , containing all the second-order derivatives as illustrated in the following drawing. The diagonal $(d \times d)$ matrices correspond to self-correlation derivatives $\frac{\partial C}{\partial w_{0,z} \partial w_{0,t}}$ and $\frac{\partial C}{\partial w_{k',z} \partial w_{k',t}}$ for $k' = 1 \dots K$. The $(d \times d)$ off-diagonals contain cross-correlation derivatives $\frac{\partial C}{\partial w_{0,z} \partial w_{k',t}}$ for $k' = 1 \dots K$.

$$\mathbf{H} = \begin{array}{ccccc} & \begin{array}{c} d \\ \boxed{\frac{\partial C}{\partial \mathbf{w}_0 \partial \mathbf{w}_0}} \end{array} & \begin{array}{c} d \\ \boxed{\frac{\partial C}{\partial \mathbf{w}_0 \partial \mathbf{w}_1}} \end{array} & \cdots & \begin{array}{c} d \\ \boxed{\frac{\partial C}{\partial \mathbf{w}_0 \partial \mathbf{w}_K}} \end{array} \\ \begin{array}{c} d \\ \boxed{\frac{\partial C}{\partial \mathbf{w}_1 \partial \mathbf{w}_0}} \end{array} & \begin{array}{c} \boxed{\frac{\partial C}{\partial \mathbf{w}_1 \partial \mathbf{w}_1}} \end{array} & \cdots & \begin{array}{c} \boxed{\frac{\partial C}{\partial \mathbf{w}_1 \partial \mathbf{w}_K}} \end{array} \\ \vdots & \vdots & \ddots & \vdots \\ \begin{array}{c} d \\ \boxed{\frac{\partial C}{\partial \mathbf{w}_K \partial \mathbf{w}_0}} \end{array} & \begin{array}{c} \boxed{\frac{\partial C}{\partial \mathbf{w}_K \partial \mathbf{w}_1}} \end{array} & \cdots & \begin{array}{c} \boxed{\frac{\partial C}{\partial \mathbf{w}_K \partial \mathbf{w}_K}} \end{array} \end{array}$$

We will derive the above three types of terms in the following.

$$1) \quad \frac{\partial C}{\partial \mathbf{w}_0 \partial \mathbf{w}_0}:$$

$$\frac{\partial C}{\partial w_{0,z} \partial w_{0,t}} = \frac{\partial}{\partial w_{0,z} \partial w_{0,t}} \sum_{p,q} \left\{ \underbrace{\frac{\mu}{4} (w_{0,p}^2 w_{0,q}^2) (\mathbf{x}_{[p]}^T \mathbf{x}_{[q]})^2}_{Z_1} + \underbrace{\frac{\mu}{2} \sum_{k'=1}^K (w_{0,p}^2 w_{k',q}^2) (\mathbf{x}_{[p]}^T \mathbf{x}_{[q]})^2}_{Z_2} \right\} \quad (9)$$

which excludes the terms in Eq. (8) that do not depend on \mathbf{w}_0 .

$$\frac{\partial Z_1}{\partial w_{0,z} \partial w_{0,t}} = \begin{cases} 2\mu w_{0,z} w_{0,t} (\mathbf{x}_{[z]}^T \mathbf{x}_{[t]})^2 & \text{if } t \neq z \\ 2\mu w_{0,z}^2 (\mathbf{x}_{[z]}^T \mathbf{x}_{[z]})^2 + \mu \sum_q w_{0,q}^2 (\mathbf{x}_{[z]}^T \mathbf{x}_{[q]})^2 & \text{if } t = z \end{cases} \quad (10)$$

$$\frac{\partial Z_2}{\partial w_{0,z} \partial w_{0,t}} = \begin{cases} 0 & \text{if } t \neq z \\ \mu \sum_q \sum_{k'=1}^K w_{k',q}^2 (\mathbf{x}_{[z]}^T \mathbf{x}_{[q]})^2 & \text{if } t = z \end{cases} \quad (11)$$

Let us denote $\mathbf{W} = \mathbf{w}_0 \mathbf{w}_0^T$ and $\mathbf{G} = \mathbf{X}^T \mathbf{X}$. Then,

$$\Rightarrow \boxed{\frac{\partial C}{\partial \mathbf{w}_0 \partial \mathbf{w}_0} = 2\mu \left[\mathbf{W} \odot \mathbf{G} \odot \mathbf{G} \right] + \text{diag}(v_1^{(0)}, \dots, v_d^{(0)})} \quad (12)$$

where $v_p^{(0)} = \frac{1}{2} \sum_q (w_{0,q}^2 + \sum_{k'=1}^K w_{k',q}^2) (\mathbf{x}_{[p]}^T \mathbf{x}_{[q]})^2 \geq 0$.

$$\mathbf{2)} \quad \frac{\partial C}{\partial \mathbf{w}_k \partial \mathbf{w}_k}:$$

$$\frac{\partial C}{\partial w_{k,z} \partial w_{k,t}} = \frac{\partial}{\partial w_{k,z} \partial w_{k,t}} \sum_{p,q} \left\{ \underbrace{\frac{\mu}{4} (w_{k,p}^2 w_{k,q}^2) (\mathbf{x}_{[p]}^T \mathbf{x}_{[q]})^2}_{K_1} + \underbrace{\frac{\mu}{2} (w_{0,p}^2 w_{k,q}^2) (\mathbf{x}_{[p]}^T \mathbf{x}_{[q]})^2}_{K_2} \right\} \quad (13)$$

which excludes the terms in Eq. (8) that do not depend on \mathbf{w}_k .

$$\frac{\partial K_1}{\partial w_{k,z} \partial w_{k,t}} = \begin{cases} 2\mu w_{k,z} w_{k,t} (\mathbf{x}_{[z]}^T \mathbf{x}_{[t]})^2 & \text{if } t \neq z \\ 2\mu w_{k,z}^2 (\mathbf{x}_{[z]}^T \mathbf{x}_{[z]})^2 + \mu \sum_q w_{k,q}^2 (\mathbf{x}_{[z]}^T \mathbf{x}_{[q]})^2 & \text{if } t = z \end{cases} \quad (14)$$

$$\frac{\partial K_2}{\partial w_{k,z} \partial w_{k,t}} = \begin{cases} 0 & \text{if } t \neq z \\ \mu \sum_p w_{0,p}^2 (\mathbf{x}_{[p]}^T \mathbf{x}_{[z]})^2 & \text{if } t = z \end{cases} \quad (15)$$

Then, for $v_p^{(k)} = \frac{1}{2} \sum_q (w_{0,q}^2 + w_{k,q}^2) (\mathbf{x}_{[p]}^T \mathbf{x}_{[q]})^2 \geq 0$ we can write

$$\Rightarrow \boxed{\frac{\partial C}{\partial \mathbf{w}_k \partial \mathbf{w}_k} = 2\mu \left[\mathbf{w}_k \mathbf{w}_k^T \odot \mathbf{G} \odot \mathbf{G} \right] + \text{diag}(v_1^{(k)}, \dots, v_d^{(k)})} \quad (16)$$

$$\mathbf{3)} \quad \frac{\partial C}{\partial \mathbf{w}_0 \partial \mathbf{w}_k}:$$

$$\frac{\partial C}{\partial w_{0,z} \partial w_{k,t}} = \frac{\partial}{\partial w_{0,z} \partial w_{k,t}} \sum_{p,q} \left\{ \underbrace{\frac{\mu}{2} (w_{0,p}^2 w_{k,q}^2) (\mathbf{x}_{[p]}^T \mathbf{x}_{[q]})^2}_T \right\} \quad (17)$$

which excludes the terms in Eq. (8) that do not depend on both \mathbf{w}_0 and \mathbf{w}_k .

$$\frac{\partial T}{\partial w_{0,z} \partial w_{k,t}} = 2\mu w_{0,z} w_{k,t} (\mathbf{x}_{[z]}^T \mathbf{x}_{[t]})^2 \Rightarrow \boxed{\frac{\partial C}{\partial \mathbf{w}_0 \partial \mathbf{w}_k} = 2\mu \left[\mathbf{w}_0 \mathbf{w}_k^T \odot \mathbf{G} \odot \mathbf{G} \right]} \quad (18)$$

Let us denote $\mathbf{D}_0 = \text{diag}(v_1^{(0)}, \dots, v_d^{(0)})$, $\mathbf{D}_k = \text{diag}(v_1^{(k)}, \dots, v_d^{(k)})$, and $\mathbf{G}^2 = \mathbf{G} \odot \mathbf{G}$. Then the Hessian matrix can be written as the sum of the following matrices.

$$\mathbf{H} = \begin{bmatrix} \mathbf{w}_0 \mathbf{w}_0^T \odot \mathbf{G}^2 & \mathbf{w}_0 \mathbf{w}_1^T \odot \mathbf{G}^2 & \mathbf{w}_0 \mathbf{w}_2^T \odot \mathbf{G}^2 & \dots & \mathbf{w}_0 \mathbf{w}_K^T \odot \mathbf{G}^2 \\ \mathbf{w}_1 \mathbf{w}_0^T \odot \mathbf{G}^2 & \mathbf{w}_1 \mathbf{w}_1^T \odot \mathbf{G}^2 & \mathbf{0} & \dots & \mathbf{0} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ \mathbf{w}_K \mathbf{w}_0^T \odot \mathbf{G}^2 & \mathbf{0} & \dots & \mathbf{0} & \mathbf{w}_K \mathbf{w}_K^T \odot \mathbf{G}^2 \end{bmatrix} + \begin{bmatrix} \mathbf{D}_0 & \mathbf{0} & \mathbf{0} & \dots & \mathbf{0} \\ \mathbf{0} & \mathbf{D}_1 & \mathbf{0} & \dots & \mathbf{0} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ \mathbf{0} & \mathbf{0} & \dots & \mathbf{0} & \mathbf{D}_K \end{bmatrix} \quad (19)$$

The second diagonal matrix, denoted \mathbf{H}_2 , is PSD as it contains non-negative entries $v_p^{(k)} = \frac{1}{2} \sum_q (w_{0,q}^2 + w_{k,q}^2) (\mathbf{x}_{[p]}^T \mathbf{x}_{[q]})^2$, $p = 1 \dots d$ for $k = 1, \dots, K$ and non-negative entries $v_p^{(0)} = \frac{1}{2} \sum_q (w_{0,q}^2 + \sum_{k'=1}^K w_{k',q}^2) (\mathbf{x}_{[p]}^T \mathbf{x}_{[q]})^2$, $p = 1 \dots d$.

We can also show that the first matrix, denoted \mathbf{H}_1 , is PSD by first decomposing it into a sum of outer products between vectors of the form

$$\mathbf{G}^2 \odot \left(\begin{bmatrix} \frac{\mathbf{w}_0}{\sqrt{(K+1)}} \\ \mathbf{w}_1 \\ \mathbf{0} \\ \vdots \\ \mathbf{0} \end{bmatrix} \left[\frac{\mathbf{w}_0}{\sqrt{(K+1)}} \mathbf{w}_1^T \mathbf{0}^T \dots \mathbf{0}^T \right] + \begin{bmatrix} \frac{\mathbf{w}_0}{\sqrt{(K+1)}} \\ \mathbf{0} \\ \mathbf{w}_2 \\ \vdots \\ \mathbf{0} \end{bmatrix} \left[\frac{\mathbf{w}_0}{\sqrt{(K+1)}} \mathbf{0}^T \mathbf{w}_2^T \dots \mathbf{0}^T \right] \right) \quad (20)$$

$$\dots + \begin{bmatrix} \frac{\mathbf{w}_0}{\sqrt{(K+1)}} \\ \mathbf{0} \\ \vdots \\ \mathbf{0} \\ \mathbf{w}_K \end{bmatrix} \left[\frac{\mathbf{w}_0}{\sqrt{(K+1)}} \mathbf{0}^T \dots \mathbf{0}^T \mathbf{w}_K^T \right] = \mathbf{G}^2 \odot (\mathbf{v}_1 \mathbf{v}_1^T + \dots + \mathbf{v}_K \mathbf{v}_K^T) = \mathbf{H}_1 \quad (21)$$

It is easy to see that $\mathbf{r}^T \mathbf{H}_1 \mathbf{r} \geq 0$ for *any* vector \mathbf{r} , since $(\mathbf{r}^T \mathbf{v}_k)(\mathbf{v}_k^T \mathbf{r}) = s_k^2 \geq 0$ for all k terms that constitute the sum.

We have shown that both \mathbf{H}_1 and \mathbf{H}_2 are PSD and since the sum of two PSD matrices is also PSD, then the Hessian matrix \mathbf{H} is also PSD, which concludes the proof for convexity. \square

A.2 Top-level classification: *Precision (seen)*, *Recall (seen)*, *Recall (unseen)* for RISK-SEN and NYT-DSTR datasets.

Figure 8 and 9 provide additional results of *Precision (seen)*, *Recall (seen)* and *Recall (unseen)* on two datasets: RISK-SEN and NYT-DSTR.

A.3 Sub-level classification: Breakdown confusion tables to compute *acc(rare)* for the three datasets.

Tables 4, 5 and 6 present confusion tables of *acc(rare)* on RISK-DOC, RISK-SEN and NYT-DSTR datasets.

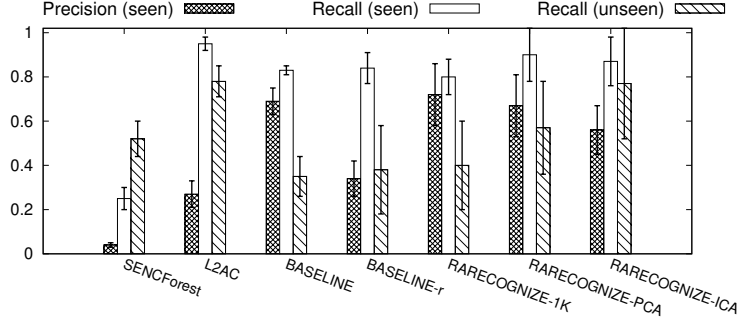


Fig. 8: Top-level classification: *Precision (seen)*, *Recall (seen)* and *Recall (unseen)* of methods on RISK-SEN dataset.

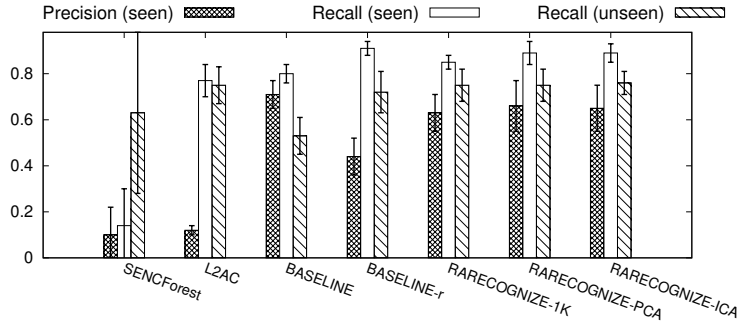


Fig. 9: Top-level classification: *Precision (seen)*, *Recall (seen)* and *Recall (unseen)* of methods on NYT-DSTR dataset.

A.4 Interpretability: Wordclouds for RARECOGNIZE in all three datasets RISK-DOC, RISK-SEN and NYT-DSTR.

Fig. 10, 11 and 12 show complete word clouds for both general and specialized classifiers for all three datasets considered in our paper.

Table 4: Sub-level classification: Breakdown confusion tables to compute $acc(rare)$ in RISK-DOC (Bolded values reflect correct predictions).

Label	\mathcal{R}_{ks}	\mathcal{R}_u	\mathcal{N}_{test}	\mathcal{R}_{ks}	\mathcal{R}_u	\mathcal{N}_{test}	\mathcal{R}_{ks}	\mathcal{R}_u	\mathcal{N}_{test}
$\hat{\mathcal{R}}_{ks}$	81.8	6.0	38.4	40.8	111.6	125.8	294.0	64.2	46.6
$\hat{\mathcal{R}}_{k's}(k' \neq k)$	8.6			225.6			14.0		
$\hat{\mathcal{R}}_u$	301.0	218.8	482.4	115.6	139.4	44.4	75.0	44.2	125.4
$\hat{\mathcal{N}}$	0.8	205.6	26.4	10.2	179.4	377.0	9.2	322.0	375.2

(a) SENCFOREST

\mathcal{R}_{ks}	\mathcal{R}_u	\mathcal{N}_{test}
294.0	70.4	46.6
14.0		
75.0	59.2	125.6
9.2	300.8	375.0

(d) BASELINE-r

(b) L2AC

\mathcal{R}_{ks}	\mathcal{R}_u	\mathcal{N}_{test}
293.0	28.0	35.2
13.2		
54.0	77.2	21.0
32.0	325.2	491.0

(e) RARECOGNIZE-1K

(c) BASELINE

\mathcal{R}_{ks}	\mathcal{R}_u	\mathcal{N}_{test}
330.4	22.8	25.0
14.8		
0.4	81.0	48.2
46.6	326.4	474.0

(f) RARECOGNIZE-PCA

\mathcal{R}_{ks}	\mathcal{R}_u	\mathcal{N}_{test}
304.0	102.8	220.2
20.0		
0.0	214.2	137.2
68.2	113.4	189.8

(g) RARECOGNIZE-ICA

Table 5: Sub-level classification: Breakdown confusion tables to compute $acc(rare)$ in RISK-SEN.

	\mathcal{R}_{ks}	\mathcal{R}_u	\mathcal{N}_{test}		\mathcal{R}_{ks}	\mathcal{R}_u	\mathcal{N}_{test}		\mathcal{R}_{ks}	\mathcal{R}_u	\mathcal{N}_{test}
$\hat{\mathcal{R}}_{ks}$	58.0	185.6	769.6		57.2	236.2	426.6		156.0	68.8	92.0
$\hat{\mathcal{R}}_{k's}(k' \neq k)$	80.0				98.8				11.0		
$\hat{\mathcal{R}}_u$	99.6	199.4	776.4		70.8	65.6	187.8		30.2	73.0	207.6
$\hat{\mathcal{N}}$	0.0	0.0	0.0		10.8	83.2	931.6		40.4	243.2	1246.4

(a) SENCFOREST	(b) L2AC	(c) BASELINE																																													
<table> <tr><th>\mathcal{R}_{ks}</th><th>\mathcal{R}_u</th><th>\mathcal{N}_{test}</th></tr> <tr><td>156.0</td><td>74.2</td><td>92.0</td></tr> <tr><td>11.0</td><td></td><td></td></tr> <tr><td>30.2</td><td>81.4</td><td>207.6</td></tr> <tr><td>40.4</td><td>229.4</td><td>1246.4</td></tr> </table>	\mathcal{R}_{ks}	\mathcal{R}_u	\mathcal{N}_{test}	156.0	74.2	92.0	11.0			30.2	81.4	207.6	40.4	229.4	1246.4	<table> <tr><th>\mathcal{R}_{ks}</th><th>\mathcal{R}_u</th><th>\mathcal{N}_{test}</th></tr> <tr><td>158.0</td><td>83.6</td><td>36.0</td></tr> <tr><td>10.6</td><td></td><td></td></tr> <tr><td>22.2</td><td>78.0</td><td>46.6</td></tr> <tr><td>46.8</td><td>223.4</td><td>1463.4</td></tr> </table>	\mathcal{R}_{ks}	\mathcal{R}_u	\mathcal{N}_{test}	158.0	83.6	36.0	10.6			22.2	78.0	46.6	46.8	223.4	1463.4	<table> <tr><th>\mathcal{R}_{ks}</th><th>\mathcal{R}_u</th><th>\mathcal{N}_{test}</th></tr> <tr><td>186.0</td><td>41.6</td><td>52.8</td></tr> <tr><td>28.2</td><td></td><td></td></tr> <tr><td>2.2</td><td>180.8</td><td>70.2</td></tr> <tr><td>21.2</td><td>161.4</td><td>1423.0</td></tr> </table>	\mathcal{R}_{ks}	\mathcal{R}_u	\mathcal{N}_{test}	186.0	41.6	52.8	28.2			2.2	180.8	70.2	21.2	161.4	1423.0
\mathcal{R}_{ks}	\mathcal{R}_u	\mathcal{N}_{test}																																													
156.0	74.2	92.0																																													
11.0																																															
30.2	81.4	207.6																																													
40.4	229.4	1246.4																																													
\mathcal{R}_{ks}	\mathcal{R}_u	\mathcal{N}_{test}																																													
158.0	83.6	36.0																																													
10.6																																															
22.2	78.0	46.6																																													
46.8	223.4	1463.4																																													
\mathcal{R}_{ks}	\mathcal{R}_u	\mathcal{N}_{test}																																													
186.0	41.6	52.8																																													
28.2																																															
2.2	180.8	70.2																																													
21.2	161.4	1423.0																																													

(d) BASELINE-r	(e) RARECOGNIZE-1K	(f) RARECOGNIZE-PCA															
	<table> <tr><th>\mathcal{R}_{ks}</th><th>\mathcal{R}_u</th><th>\mathcal{N}_{test}</th></tr> <tr><td>194.6</td><td>45.4</td><td>58.8</td></tr> <tr><td>20.2</td><td></td><td></td></tr> <tr><td>0.0</td><td>198.0</td><td>74.6</td></tr> <tr><td>20.4</td><td>150.8</td><td>1459.2</td></tr> </table>	\mathcal{R}_{ks}	\mathcal{R}_u	\mathcal{N}_{test}	194.6	45.4	58.8	20.2			0.0	198.0	74.6	20.4	150.8	1459.2	
\mathcal{R}_{ks}	\mathcal{R}_u	\mathcal{N}_{test}															
194.6	45.4	58.8															
20.2																	
0.0	198.0	74.6															
20.4	150.8	1459.2															

(g) RARECOGNIZE-ICA

Table 6: Sub-level classification: Breakdown confusion tables to compute $acc(rare)$ in NYT-DSTR.

	\mathcal{R}_{ks}	\mathcal{R}_u	\mathcal{N}_{test}		\mathcal{R}_{ks}	\mathcal{R}_u	\mathcal{N}_{test}		\mathcal{R}_{ks}	\mathcal{R}_u	\mathcal{N}_{test}
$\hat{\mathcal{R}}_{ks}$	0.2	1.0	0.0		3.0	420.8	1246.6		158.4	132.2	175.2
$\hat{\mathcal{R}}_{k's}(k' \neq k)$	0				175.6				12.6		
\mathcal{R}_u	291.0	330.4	2112.4		47.0	77.2	337.6		64.0	234.4	774.0
\mathcal{N}	0.2	340.0	0.6		65.8	163.4	528.8		56.4	294.8	1163.8

(a) SENCFOREST

\mathcal{R}_{ks}	\mathcal{R}_u	\mathcal{N}_{test}
158.4	135.8	175.2
12.6		
64.0	243.2	774.0
56.4	282.4	1163.8

(b) L2AC

\mathcal{R}_{ks}	\mathcal{R}_u	\mathcal{N}_{test}
176.2	251.0	69.0
16.4		
56.2	260.2	77.4
42.6	150.2	1966.6

(c) BASELINE

\mathcal{R}_{ks}	\mathcal{R}_u	\mathcal{N}_{test}
221.4	57.2	85.4
44.2		
0.0	400.0	129.0
25.8	204.2	1898.6

(d) BASELINE-r

(e) RARECOGNIZE-1K

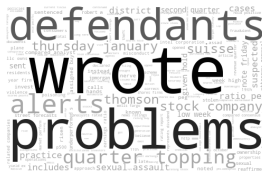
\mathcal{R}_{ks}	\mathcal{R}_u	\mathcal{N}_{test}
213.2	62.4	296.2
52.6		
0.0	393.6	0.4
25.6	205.4	1816.4

(f) RARECOGNIZE-PCA

(g) RARECOGNIZE-ICA



Fig. 10: Interpretability: Word clouds representing learned weights by general and all specialized classifiers on NYT-DSTR dataset.



(a) General risk



(b) Climate change



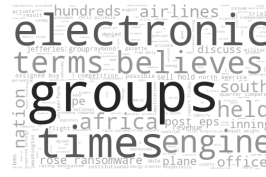
(c) Cyber-attack



(d) Data leaks



(e) Drug abuse



(f) Engine failure



(g) Financial fraud



(h) Gun violation



(i) Low stock rating



(j) Military attack



(k) Misleading statement



(l) Money laundering



(m) Negative earning



(n) Sexual assault



(o) Spying



(p) Trade war

Fig. 11: Interpretability: Word clouds representing learned weights by general and all specialized classifiers on RISK-DOC dataset.



(a) General risk



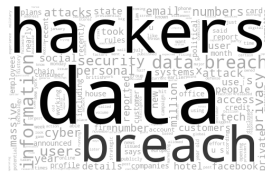
(b) Bankruptcy



(c) Corruption



(d) Counterfeiting



(e) Cyber-privacy



(f) Environment



(g) Fraud False Claims



(h) Labor



(i) Money laundering

Fig. 12: Interpretability: Word clouds representing learned weights by general and all specialized classifiers on RISK-SEN dataset.