

LIPEX – LOCALLY INTERPRETABLE PROBABILISTIC EXPLANATIONS – TO LOOK BEYOND THE TRUE CLASS

Anonymous authors

Paper under double-blind review

ABSTRACT

In this work, we instantiate a novel perturbation-based multi-class explanation framework, LIPEX (**L**ocally **I**nterpretable **P**robabilistic **E**xplanation). We demonstrate that LIPEX not only locally replicates the probability distributions output by the widely used complex classification models but also provides insight into how every feature deemed to be important affects the prediction probability for each of the possible classes. We achieve this by defining the explanation as a matrix obtained via regression with respect to the Hellinger distance in the space of probability distributions. Ablation tests on text and image data, show that LIPEX-guided removal of important features from the data causes more change in predictions for the underlying model than similar tests on other saliency-based or feature importance-based XAI methods. It is also shown that compared to LIME, LIPEX is much more data efficient in terms of the number of perturbations needed for reliable evaluation of the explanation.

1 INTRODUCTION

Recent momentum in deep learning research has made interpreting models with complex architectures very important. In a wide range of areas where neural nets have made a successful foray, the method of “Explainable A.I.” (XAI) has also found an important use to help understand the functioning of these novel predictors - like in climate science (Labe & Barnes, 2021), for solving partial differential equation (Linial et al., 2023), in high-energy physics (Neubauer & Roy, 2022), information retrieval (Lyu & Anand, 2023), in legal A.I. (Collenette et al., 2023), etc. Most often, it has been observed that models with complex architectures give better accuracy compared to a simple model. So, the core puzzle that XAI can be seen to solve is to give a highly accurate local replication of a complex predictor’s behaviour by a simple model over humanly interpretable components of the data (Ribeiro et al., 2016). Towards achieving this, multiple different XAI methods have been proposed in the recent times, e.g., LIME (Ribeiro et al., 2016), SHAP (Lundberg & Lee, 2017), Decision-Set (Lakkaraju et al., 2016), Anchor (Ribeiro et al., 2018), Smooth-GRAD (Smilkov et al., 2017b), Poly-CAM (Englebert et al., 2022), Extremal Perturbations (Fong et al., 2019), Saliency Maps (Simonyan et al., 2014), etc.

One major motivation for explainability is debugging a model (Casillas et al., 2003; Dapaah & Grabowski, 2016). Towards this, an end user is interested not only in understanding the explanation provided for the predicted class at a particular data point but also in the influence of different features for all possible class likelihoods estimated by a classifier. The full spectrum of feature influence on each class at a particular data point can help to understand how well the model has been trained to discriminate a particular class from the rest. However, existing explanation frameworks do not provide any clue on the aforementioned issue. To this end, we propose an explainability framework that can explain a classifier’s output prediction beyond the true class.

To obtain an explanation around a data point, a local explanation algorithm like LIME (Garreau & Luxburg, 2020) creates perturbations around it, each perturbation being represented as a Boolean vector. LIME includes a feature selection method to decide a set of important features for each class (like Algorithm A) among which the perturbations are considered. Then, an explanation vector for the complex model’s prediction on the input data is obtained by solving a penalized linear regression over these perturbations and the complex classifier’s predictions on the data corresponding to the perturbations. We posit that it is not entirely convincing that LIME attempts to regress over bounded labels, i.e., probabilities, using an unbounded function (i.e., a linear function) and that this would need to be called separately for each class. Further, even if by repeated calls on each possible class

we obtain an explanation for of the classes, there is no guarantee that by these repeated evaluations, the importance of any particular feature would be knowable for every class.

In this work, we attempt to remedy these problems by proposing a single unified framework that applies to both text and images, which we will show in experiments to be better than various XAI methods for both text and images. In a C -class classification task, for any data s which is represented as \mathbf{z}_s in some f_s dimensional feature space, we shall seek explanations that map into C -class probability space as,

$$\mathbb{R}^{f_s} \ni \mathbf{z}_s \mapsto \text{Soft-Max} \circ \mathbf{W}\mathbf{z}_s \tag{1}$$

We call the $\mathbf{W} \in \mathbb{R}^{C \times f_s}$ as the “explanation matrix” – which can be obtained by minimizing some valid distance function (like Hellinger’s distance) between distributions obtained as above and the probability distribution over classes that the complex model has been trained to map any input. Thus, we instantiate this novel mechanism for XAI, namely LIPEX.

Note how the matrix \mathbf{W} in Equation 1 simultaneously gives for every feature a numerical measure of its importance for each possible class. We posit that it is important that in any explanation, it should be evident that most features deemed to be important for the predicted class are not so for the other classes - an idea that was recently formalized in Gupta & Arora (2020); Gupta et al. (2022) for the specific case of saliency maps. In our method, this property turns out to be emergent as a consequence of the more principled definition of explanation that we start from.

Figure 1 shows an example of our matrix explanation obtained for a text document. We observe how the explanation matrix is obtained for a specific document over a set of feature words. Note that for the first row (the top predicted class), the top 5 feature words detected for this instance ([feel, valued, joy, treasures, incredibly]) are *distinctly different* from the top features detected for the class in the second row, the one with the second highest probability predicted by the classifier. More examples like this can be found in Appendix D.5 and Appendix D.6 particularly focuses on examples where the predicted class and the true class are different. It is observed that there always arises a natural discrimination between features important for the different classes.

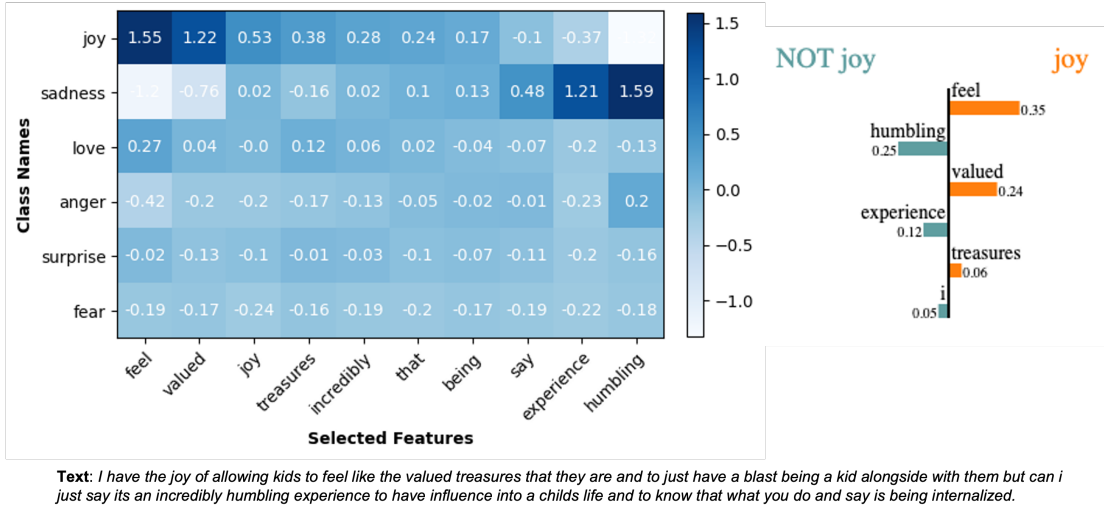


Figure 1: Example of comparison of explanation matrix obtained by LIPEX and the bar chart obtained by LIME on a text data from the Emotion dataset. For the LIPEX matrix, the class names on the left side are arranged in descending order of the predicted class probabilities. Examples of LIPEX explanation on image data is provided in Appendix D.3.

In the following, we summarize our contributions towards formalizing this idea of matrix-based explanations that match the distribution over classes predicted by a complex model.

Novel Explanation Framework In Section 3, we formally state our explanation approach, which can extract the relative importance of a set of features for every class under consideration beyond the true class. In the following tests, we demonstrate how such a multi-class explanation framework can be more useful for model understanding compared to existing state-of-the-art XAI methods.

(Test 1) Evidence of LIPEX Replicating the Complex Model’s Predicted Distribution Over Classes In Figure 2, we calculate the Total Variation (TV) distance of the output distribution of the obtained LIPEX explainer and the distribution output by the complex model for the same data and we show that over hundreds of randomly chosen test instances, the distance is overwhelmingly near 0. We show that this very necessary property holds over multiple models over text as well as images.

(Test 2) Sanity Check of LIPEX’s Sensitivity to Model Distortion In Figure 3 we distort a well-trained complex model by adding mean zero noise to the parameters in the last layer and measure how upon increasing the noise variance, the output probability distribution moves away in TV distance from the original prediction. We show that when LIPEX is implemented on the distorted models, our explainer’s predicted distribution moves away from its original value almost identically. This sanity check is inspired by the arguments in Adebayo et al. (2018).

The above two tests give robust evidence that, indeed LIPEX is an accurate local approximator of the complex model while being dramatically simpler than the black-box predictor. To the best of our knowledge, such a strong model replication property is not known to be true for even the saliency methods, which can in-principle be called on different classes separately to get the relative importance among the different pixels for each class - albeit separately.

(Test 3) Evidence of Changes in the Complex Model’s Prediction Under LIPEX Guided Data Distortion In Table 1 and 2, we devise an ablation study guided by the “faithfulness” criteria as outlined in Atanasova et al. (2020). We establish that the top features detected by LIPEX are more important for the complex model than those detected by other XAI methods. We show this by demonstrating that when the top features are removed from the data and inference is done on this damaged data, then the new predicted class differs more from the original prediction when the removal is guided by what LIPEX deemed to be important than other XAI methods.

(Test 4) Evidence of LIPEX Replicating the Complex Model’s Class Prediction Under LIPEX Guided Data Distortion In Table 3, we demonstrate that for an overwhelming majority of data, upon removing their features deemed important by LIPEX, the new class predicted by the complex predictor is reproduced by the LIPEX model when presented with the same distorted data.

(Test 5) Stability of LIPEX to Choosing Less and Only Near-Truth Perturbations In Figure 4, we demonstrate experiments that the features picked out by the LIPEX matrix are largely stable when the matrix is derived using only a few perturbation instances. We also show that this property is not true for LIME in the models we consider. Thus LIPEX is demonstrably more data efficient.

To put the above in context we recall that estimates were given in Agarwal et al. (2021) for how many perturbations around the true data are sufficient for LIME to produce reliable results - and this experiment of ours can be seen to corroborate that. Also, we recall that in works like Slack et al. (2020) it was pointed out that LIME’s reliance on perturbations far from the true data creates a vulnerability that can be exploited to create adversarial attacks.

Note that we have restricted our attention to “intrinsic evaluations” of explanations, i.e., we only use calls to the model as a black-box for deciding whether the explanations obtained are meaningful as opposed to looking for external human evaluation. Both text and image data were used to evaluate our proposed approach. For text-based experiments we used 20Newsgroup¹ and Emotion² datasets. For image-based experiments, we have used the Imagenette³ dataset with segments detected by “segment anything”⁴.

Among the above experiments, LIPEX was compared against a wide range of state-of-the-art explanation methods for both text and image data, i.e., LIME (Ribeiro et al., 2016), Guided Backpropagation (Springenberg et al., 2014), Vanilla Gradients (Erhan et al., 2009), Integrated Gradients (Sundararajan et al., 2017), DeepLift (Shrikumar et al., 2016), Occlusion (Zeiler & Fergus, 2014), XRAI (Kapishnikov et al., 2019), GradCAM (Selvaraju et al., 2017), GuidedIG (Kapishnikov et al., 2021), BlurIG (Xu et al., 2020) and SmoothGrad (Smilkov et al., 2017a).

¹<http://qwone.com/~jason/20Newsgroups/>

²<https://huggingface.co/datasets/dair-ai/emotion>

³<https://github.com/fastai/imagenette>

⁴<https://segment-anything.com/>

Organization In Section 2 we briefly overview related works in XAI. In Section 3 we give the precise loss function formalism for obtaining our explanation matrix, and in Section 4 all the tests will be given - comparing the relative benefits to other XAI methods. We conclude in Section 5. Appendices contain various details such as the precise pseudocode used in Section 4 (in Appendix C), the hyperparameter settings (in Appendix B), and further experimental data is given in Appendix D.

2 RELATED WORK

The work in Letham et al. (2015) is one of the first works that attempted to develop a classifier using rules and Bayesian analysis. In Ribeiro et al. (2016) a first attempt was made to describe explainability formally. The explanation can be made through an external explainer module, or a model can also be attempted to be made inherently explainable (Chattopadhyay et al., 2023). Post-hoc explainer strategy, as is the focus here, can be of different types, like (a) Ribeiro et al. (2016); Lundberg & Lee (2017) estimate feature importance for predicting a particular output, (b) counterfactual explanations (Wachter et al. (2017); Ustun et al. (2019); Rawal & Lakkaraju (2020)) determine if a feature x was present in the input, then would the model have predicted output y , (c) contrastive approaches (Jacovi et al., 2021) describe why an ML model has predicted a particular output instead of another, or (d) Weinberger et al. (2023) and Crabbé & van der Schaar (2022) have recently proposed new XAI methods tuned to the case of unsupervised learning. In this work, we specifically focus on feature importance-based explanation techniques.

Feature Importance-based Explanations The study in Ribeiro et al. (2016) initiated the LIME framework which we reviewed in Section 1 as our primary point of motivation. Similarly, the work in Lundberg & Lee (2017) used a statistical sampling approach (“SHAP”) to explain a classifier model in terms of human interpretable features. Lakkaraju et al. (2016) proposed a decision set-based approach to train a classifier that can be interpretable and accurate simultaneously - where a set of independent if-then rules defines a decision set. Ribeiro et al. (2018) proposed an anchor-based approach for explanation - where anchors were defined as a set of sufficient conditions for a particular local prediction.

Evaluation is a critical component in any explanation framework. The study in Doshi-Velez & Kim (2017) described important characteristics for the evaluation of explanation approaches. Evaluation criteria for explanations can broadly be categorized into two types, (a) criteria which measure how well the explainer module is able to mimic the original classifier and (b) criteria which measure the trustworthiness of the features provided by the explainer module, like the work in Qi et al. (2019) demonstrated the change in the prediction probability of a classifier with the removal of top K features predicted by a saliency map explainer.

We note that in this work our tests done in Section 4 encompass both the above kinds of criteria.

Lastly, we note that in Sokol & Flach (2020) a tree based explanation was attempted which could directly work in the multiclass setting but to be able to compete LIME their method’s computation cost can need to scale with the number of segments in an image. Also, in sharp contrast to our LIPEX proposal, it does not have the critical ability to explain/reproduce the predicted distribution of the given complex model.

3 OUR SETUP

Let $\mathcal{C} \in \{1, 2, 3, \dots\}$ be the number of classes in the classification setup. Given any two probability vectors $\mathbf{p}, \mathbf{q} \in [0, 1]^{\mathcal{C}}$, $\sum_{i=1}^{\mathcal{C}} p_i = 1 = \sum_{i=1}^{\mathcal{C}} q_i$, we succinctly represent \mathbf{p}, \mathbf{q} as being members of the simplex in \mathcal{C} -dimensions as $\mathbf{p}, \mathbf{q} \in \Delta^{\mathcal{C}}$.

Classifier Setup We aim to explain a classifier which can be described as a neural network $\mathcal{N}_{\mathbf{w}}$ (parameterized by weight \mathbf{w}) composed with a layer of soft-max so that the output of the composition is a probability distribution over the \mathcal{C} -classes. Thus we define the composed mapping,

$$\mathbf{f}_{\mathbf{w}} : \mathbb{R}^d \rightarrow \Delta^{\mathcal{C}}, \mathbf{x} \mapsto \text{Soft-Max} \circ \mathcal{N}_{\mathbf{w}}(\mathbf{x}) \quad (2)$$

This composed function $\mathbf{f}_{\mathbf{w}}$ in Equation 2 commonly would have been trained via the cross-entropy loss on a \mathcal{C} class labeled data - and we assume only black-box access to it.

The Feature Space for Explanations For a specific data s (e.g., a piece of text), we denote the number of unique features (e.g., words) as $|s|$ and assume that there is a selected ‘feature space’ with f_s features. Suppose special subsets of them, say $\mathcal{S}(s)$ and $\mathcal{S}_f(s)$ have been chosen and there is a map, say Select which does the feature selection for each of its domain points as per say Algorithm A.

$$\mathcal{S}(s) \subset \mathbb{R}^{|s|}, \mathcal{S}_f(s) \subset \mathbb{R}^{f_s} \ \& \ \text{Select} : \mathcal{S}(s) \rightarrow \mathcal{S}_f(s) \quad (3)$$

The Local Explanation Matrix We explain f_w ’s behaviour around s by a ‘pseudo-linear model’, $g_{s,w}$ which is defined as,

$$g_{s,w} : \mathcal{S}_f(s) \rightarrow \Delta^C, \mathbf{z}' \mapsto \text{Soft-Max} \circ \mathbf{W}\mathbf{z}' \quad (4)$$

with $\mathbf{W} \in \mathbb{R}^{C \times f_s}$ being the ‘‘explanation matrix’’.

In the LIME setup (as well as in LIPEX), $\mathcal{S}(s) \subseteq \{0, 1\}^{|s|}$ i.e. Boolean vectors are used to represent random ways of dropping one or more of the (unique) words for text data and pixels for image data. Hence, in such setups, the original input instance is represented as an all-ones vector, $\mathbf{1}_s \in \mathbb{R}^{|s|}$.

We assume that there is a pre-chosen function (say T_s) that maps ‘‘perturbations’’ of the data contained in the set $\mathcal{S}(s)$ to some d -dimensional embedding (like the BERT embeddings) which can be input to the original prediction model (Equation 2).

$$T_s : \mathcal{S}(s) \rightarrow \mathbb{R}^d \quad (5)$$

Note that, LIME seeks explanations using a linear function which would map the \mathbf{z}' (as in Equation 4) to a real number which is a priori unbounded in sharp contrast to the explainer $g_{s,w}$ defined in Equation 4. Also note that the input dimensions d of f_w and f_s for $g_{s,w}$ could be very different and dependent on s and typically, $f_s \ll d$. Eg., in standard LIME implementations for a classifier one often chooses $f_s = 6$ important features of the text s .

The space of all probability distributions admits various natural metrics and Hellinger distance has previously been used for feature selection in classification (Fu et al., 2020). Hellinger distance between two discrete distributions \mathbf{p}, \mathbf{q} (on a set of \mathcal{C} possible classes) is given as,

$$H(\mathbf{p}, \mathbf{q}) := \frac{1}{\sqrt{2}} \cdot \sqrt{\sum_{c \in \mathcal{C}} (\sqrt{\mathbf{p}(c)} - \sqrt{\mathbf{q}(c)})^2}$$

Apart from being an intuitive symmetric measure, squared Hellinger distance also offers other attractive features of being sub-additive, smaller than half of the KL divergence and always being within a quadratic factor of the Total Variation (TV) distance. (Canonne, 2020)⁵

Let $\tilde{\mathcal{S}}(s) \subset \mathcal{S}(s) (\subseteq \{0, 1\}^{|s|})$ be a randomly sampled set of perturbations to be used for training. Passing it through the Select map (Equation 3) we obtain $\tilde{\mathcal{S}}_f(s) \subset \mathcal{S}_f(s) (\subset \mathbb{R}^{f_s})$ which are the feature representations of the perturbations. We posit that the outputs of the Select map would determine what the explainer $g_{s,w}$ in Equation 4 acts on. Further noting that the output of the embedding map T_s in Equation 5 determines what the true predictor f_w gets as input, we consider the following empirical risk function corresponding to a distance function π in $\mathbb{R}^{|s|}$,

$$\hat{\mathcal{L}}_H(g_{s,w}, \tilde{\mathcal{S}}(s)) = \frac{1}{|\tilde{\mathcal{S}}(s)|} \sum_{\mathbf{x} \in \tilde{\mathcal{S}}(s)} \pi(\mathbf{1}_s, \mathbf{x}) \cdot H^2(g_{s,w} \circ \text{Select}(\mathbf{x}), f_w \circ T_s(\mathbf{x})) + \frac{\lambda}{2} \cdot \|\mathbf{W}\|_F^2 \quad (6)$$

where $\mathbf{1}_s$, the all-ones vector in $\mathbb{R}^{|s|}$. We choose $\pi(\mathbf{1}_s, \mathbf{x}) = 1 - \frac{\mathbf{x}^\top \mathbf{1}_s}{\|\mathbf{x}\| \cdot \|\mathbf{1}_s\|}$ for all our experiments. It is immediately interpretable that Equation 6 takes a π -weighted empirical average of the Hellinger distance squared between the true distribution over classes predicted by the complex classifier f_w and the distribution predicted by the explainer $g_{s,w}$ while the λ -term penalizes for using high weight explainers and hence promotes simplicity of $g_{s,w}$.

⁵Our experiments were tried with TV and they underperformed compared to the squared Hellinger metric.

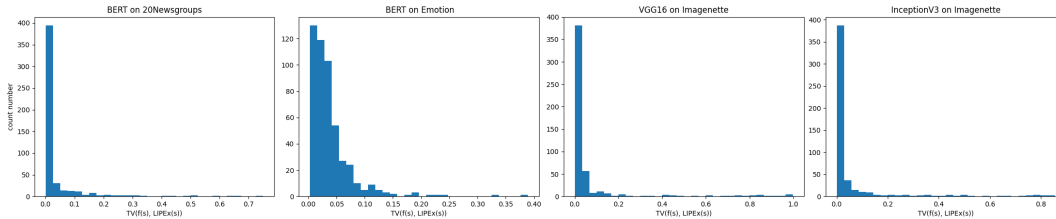


Figure 2: Histogram Statistics of TV Distance on Probability Distribution of Classes Between Classifier and LIPEX.

Intuition for Good LIPEX Minima for Text Classifiers Being a ReLU Net For intuition, consider explaining classification predictions by a ReLU net on a text data s with $|s|$ unique words. Further suppose the classifier has been trained to accept $d(\geq |s|)$ word length texts in their TF-IDF representation. Thus the T_s map (Equation 5) that lifts the perturbations of s to the input space of the complex classifier can be imagined as a tall matrix of dimensions $d \times |s|$ whose top $|s| \times |s|$ block is a diagonal matrix giving the TF-IDF values for the words in this text and the rest of the matrix being zeros. Also, we note that the Select function can be imagined as a linear projection of the Boolean-represented perturbations into the subset of important features.

Further, any ReLU neural net is a continuous piecewise linear function Arora et al. (2018). Hence, except at the measure zero set of non-differentiable points, the function \mathcal{N}_w (Equation 2) is locally a linear function. Thus, for almost every input $\mathbf{z} \in \mathbb{R}^d$ there exists a (possibly small) neighbourhood of it where $\mathcal{N}_w = \mathbf{W}_{\text{net}}$ for some matrix $\mathbf{W}_{\text{net}} \in \mathbb{R}^{C \times d}$. It would be natural to expect that most true texts are not at the non-differentiable points of the net’s domain and that T_s maps small perturbations of the data into a small neighbourhood. Hence, for many perturbations $\mathbf{x} \in \mathbb{R}^{|s|}$, $\mathbf{f}_w(T_s(\mathbf{x}))$, as it occurs in the loss in Equation 6, is a Soft-Max of a linear transformation (composition of the net and the T_s map) of \mathbf{x} . Recall that this is exactly the functional form of the explainer $\mathbf{g}_{s,W}$ (Equation 4) given that the Select function can be represented as a linear map! Thus, we see that there is a very definitive motivation for this loss function to yield good locally linear explanations for ReLU nets classifying text.

4 RESULTS

At the very outset, we note the following salient points about our setup. *Firstly*, that for any data s (say a piece of text or an image), when implementing LIPEX on it, we generate a set of 1000 perturbations of input instances. Then we chose features by taking a union set over the top-3 features of each possible class, which was returned by the “forward feature selection” method (reproduced in Algorithm A) called on the above perturbation data set. We recall that this feature selection algorithm is standard in LIME implementations⁶. Suppose this union has f_s features - then for all computations to follow for s we always stick with these f_s features for LIPEX (and also always call LIME on f_s number of features in comparison experiments). *Secondly*, we note that for the matrix returned by LIPEX (i.e. \mathbf{W} in Equation 4) we shall define its “top- k ” features as the features/columns of the matrix which give the k -highest entries by absolute value for the predicted class of that data.

Reproducing the Distribution over Classes of the Complex Classifier A key motivation for introducing the LIPEX framework was the need for the explanation framework to produce class distributions closely resembling those of the original classifier. Therefore, our initial emphasis is on investigating how much in Total Variation (TV) distance, the distribution over classes predicted by the obtained explainer is away from the one predicted for the same data by the complex model needing explanations. In Figure 2, we show the statistics of this TV distance for experiments on both text (i.e. BERT on 20NewsGroups and BERT on Emotion) and image data (i.e. VGG16 and InceptionV3 on Imagenette). Figure 2 clearly shows that the distribution is highly skewed towards 0 over five hundred randomly sampled data over multiple modalities and state-of-the-art models. Note that the LIPEX loss (Equation 6) never directly optimized for the TV gap to be small and hence we posit that this is a strong test of performance that LIPEX passes.

⁶<https://github.com/marcotcr/lime>

LIPEX Tracks Distortions of The Complex Model’s Output Distribution This sanity check experiment is inspired by the studies in Adebayo et al. (2018). Here, we add mean-zero Gaussian noise to the trained complex model’s last-layer weights and keep dialling up the noise variance till the model’s accuracy is heavily damaged. At each noise level we compute the average over randomly sampled data, of the Total Variation distance between the output distribution of the damaged model and its original value and the same for the LIPEX’s distribution for that model at respective inputs. We do text experiments with BERT on the Emotion dataset and image experiments with VGG16

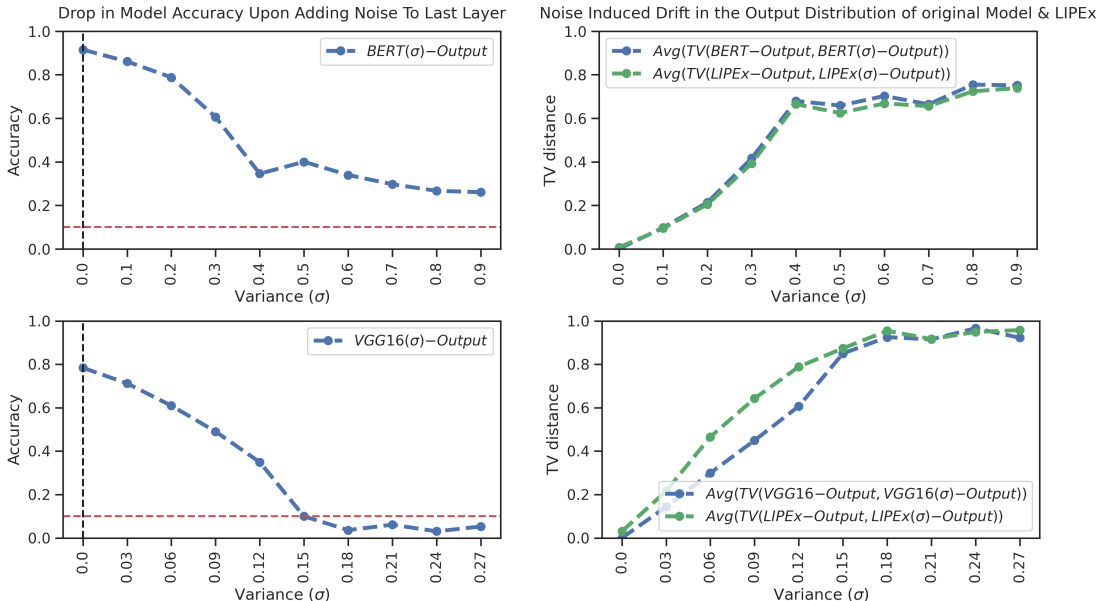


Figure 3: In the left image, we see how the model accuracy drops upon adding noise to the last layer weights and bias, a mean 0 Gaussian noise at different variances. This demonstrates that the maximum added noise is sufficient to distort the model highly. On right, we see how the data averaged TV distance between the output distribution of LIPEX for the data at the original complex model and the noise-distorted complex model tracks the same change in the complex model’s output.

on the Imagenette dataset. In Figure 3, for any specific data, LIPEX-Output is the LIPEX’s output distribution for the original model, LIPEX(σ)-Output is the LIPEX’s output distribution for the distorted model at noise variance σ . BERT-Output, BERT(σ)-Output, VGG16-Output and VGG16(σ)-Output are defined similarly.

The right column of plots in Figure 3 demonstrates that as the model distorts, LIPEX’s output moves away from the original in a remarkably identical fashion as the distorted model’s output changes with respect to its original value.

Importance of Top-K Features Detected by LIPEX A test of the correctness of determining any set of features to be important by an explanation method is that upon their removal from the original data and on presenting this modified/damaged input to the complex classification model it should produce a new predicted class than originally. We implement this test with text data in Table 1 and with image data in Table 2. We demonstrate that when the top features detected by LIPEX are removed from the data, the original model’s predicted class changes substantially more than when the same is measured for many other XAI methods for the predicted class - and the amount of change is proportional to the number of top features removed.⁷

LIPEX Reproduces the Complex Model’s Class Predictions Under LIPEX Guided Data Damage We posit that for a multi-class explainer as LIPEX, it is a very desirable sanity check that it should reproduce the underlying model’s (new) predicted classes on the input when its top features are

⁷We use the code in Atanasova et al. (2020) to implement the gradient-based methods in Table 1, and the package <https://github.com/PAIR-code/saliency> to implement the saliency methods in Table 2

Model & Dataset	Top-K	LIPEX	LIME	GuidedBack	Saliency	InputXGrad	Deeplift	Occlusion
BERT 20NewsGroups	K=1	0.781	0.777	0.387	0.387	0.38	0.38	0.45
	K=2	0.857	0.841	0.477	0.477	0.48	0.48	0.543
	K=3	0.897	0.856	0.517	0.517	0.52	0.52	0.59
	K=4	0.909	0.881	0.517	0.517	0.523	0.523	0.627
	K=5	0.908	0.912	0.553	0.553	0.57	0.57	0.653
BERT Emotion	K=1	0.657	0.653	0.597	0.597	0.6	0.6	0.65
	K=2	0.74	0.697	0.61	0.61	0.62	0.63	0.66
	K=3	0.73	0.647	0.637	0.637	0.637	0.653	0.697
	K=4	0.73	0.64	0.623	0.623	0.633	0.643	0.697
	K=5	0.793	0.65	0.63	0.63	0.637	0.64	0.693

Table 1: Here, features refer to words. Upon removing top-K words detected by each of the XAI methods and doing re-prediction, we report the fraction of data on which the predicted class changes. We see that the words removed by LIPEX guidance more significantly impact the model’s prediction than when guided by the other XAI methods. The complete experimental data with standard deviations can be seen in Table 6 in the appendix.

Model & Dataset	Top-K	LIPEX	LIME	XRAI	GradCAM	GuidedIG	BlurIG	VanillaGrad	SmoothGrad	IG
VGG16 Imagenette	K= 2	0.763	0.74	0.713	0.69	0.717	0.713	0.68	0.747	0.703
	K= 3	0.82	0.78	0.77	0.763	0.75	0.787	0.753	0.817	0.747
	K= 4	0.867	0.793	0.793	0.79	0.793	0.807	0.807	0.843	0.773
InceptionV3 Imagenette	K= 2	0.673	0.63	0.693	0.653	0.663	0.647	0.657	0.65	0.637
	K= 3	0.753	0.713	0.7	0.697	0.67	0.703	0.653	0.683	0.707
	K= 4	0.773	0.767	0.74	0.72	0.713	0.717	0.72	0.74	0.733

Table 2: Here, features refer to image segments which were gotten by Segment Anything. LIPEX and LIME can be used to directly get a weight for each segment while for the saliency-based methods a segment’s importance is determined as the sum of the weights assigned to its pixels. In the table above we can see that the fraction of data on which label prediction changes under deletion of top features detected by LIPEX is consistently higher than for other XAI methods. The complete experimental data with the standard deviation can be found in the Table 7 in the appendix.

removed. In Table 3, we show with text as well as image data, that this class prediction matching holds for the LIPEX explainer for an overwhelming majority of data.

Evidence for Data Efficiency of LIPEX as Compared to LIME Since LIPEX and LIME, both are perturbation based methods, a natural question arises if LIPEX is more data-efficient, or in other words can its top features detected be stable if only a few perturbations close to the true data are allowed. In this test, we show that not only is this true, but also that (a) LIPEX’s top features can at times even remain largely invariant to reducing the perturbations and also that (b) the difference with respect to LIME in the list of top features detected, is maintained when the allowed set of perturbations are increasingly constrained to be few and near the true data. Our comparison method is specified precisely as Algorithm C in the Appendix and we sketch it here as follows.

When in the setting with unrestricted perturbations, we infer two lists of top features, one from the row of the predicted class of the matrix (i.e. W) returned by LIPEX and another from LIME’s weight vector for the same class — say LIPEX-List- s and LIME-List- s respectively. Next, we parameterize the restriction on the allowed perturbations by the maximum angle δ that any Boolean vector representing the perturbation is allowed to subtend with respect to the all-ones vector that represents the input data.

We use the default set of perturbations in a LIME implementation as a baseline ⁸ and at different δ , we use only the δ -restricted subset of the perturbations to compute (for the model predicted class) the top features returned by the LIPEX matrix and the LIME, say δ -LIPEX-List- s and δ -LIME-List- s respectively. For quantifying the dissimilarities between these lists of top features measured by the two methods, we compute the following Jaccard indices and average the results on 100 randomly chosen instances.

⁸In the LIME code, they choose perturbations of 5000 for text data and 1000 for image data.

Model & Dataset	Modality	Top1	Top2	Top3	Top4	Top5
BERT (20NewsGroups)	Text	0.90 (± 0.041)	0.85 (± 0.024)	0.79 (± 0.039)	0.71 (± 0.033)	0.70 (± 0.005)
BERT (Emotion)	Text	0.89 (± 0.022)	0.84 (± 0.025)	0.84 (± 0.017)	0.82 (± 0.037)	0.74 (± 0.033)
VGG16 (Imagenette)	Image	0.80 (± 0.046)	0.73 (± 0.034)	0.73 (± 0.034)	0.73 (± 0.025)	0.70 (± 0.075)
InceptionV3 (Imagenette)	Image	0.90 (± 0.051)	0.78 (± 0.044)	0.75 (± 0.015)	0.74 (± 0.013)	0.69 (± 0.035)

Table 3: In this table, for each model and data combination, we give the fraction of data (over 100 random samples) over which the new class predicted by the complex model matches the new prediction by the LIPEX for the same model, upon removing from the data its top features as determined by LIPEX. We can observe that post this distortion on the data, the class labels from the complex model match those from the simple explainer for a significant majority of the instances.

$$J_{s,\delta,LIME} := \frac{|\delta\text{-LIME-List-s} \cap \text{LIME-List-s}|}{|\delta\text{-LIME-List-s} \cup \text{LIME-List-s}|}, J_{s,\delta,LIPEX} := \frac{|\delta\text{-LIPEX-List-s} \cap \text{LIPEX-List-s}|}{|\delta\text{-LIPEX-List-s} \cup \text{LIPEX-List-s}|}$$

$$J_{s,\delta\text{-LIPEX-vs-LIME}} := \frac{|\delta\text{-LIPEX-List-s} \cap \text{LIME-List-s}|}{|\delta\text{-LIPEX-List-s} \cup \text{LIME-List-s}|}$$

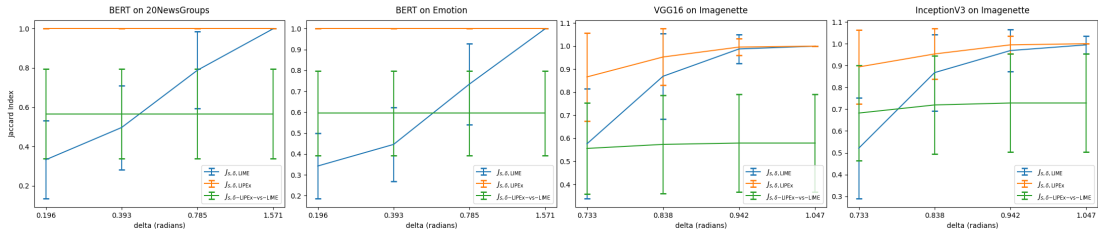


Figure 4: It can be observed that $J_{s,\delta,LIPEX}$ (the orange line) is very stable compared to that of LIME despite the allowed perturbations being made constrained. The difference in LIPEX’s features w.r.t LIME is also maintained. The number of points considered at different δ is given in Table 5 in the appendix.

From Figure 4, we can infer that at all levels of constraint on the data at least 50% of the top features detected by our LIPEX are different from LIME. *Secondly*, $J_{s,\delta,LIME}$ (averaged) rapidly falls as the number of training data allowed near the input instance is decreased. Thus its vividly revealed that the features detected by LIME are significantly influenced by those perturbations that are very far from the true text.

Lastly, and most interestingly, we note that the curve for $J_{s,\delta,LIPEX}$ (the top orange line) is very stable to using only a few perturbations which subtend a low angle with the true text. Hence the top features detected by our explanation matrix are not only important (as demonstrated in the previous two experiments) – but can also be computed very data efficiently.

5 CONCLUSION

In this work, we proposed a novel explainability framework, LIPEX, that when implemented in a classification setting, in a single training gives a weight assignment for all the possible classes for an input with respect to a chosen set of features. Unlike other XAI methods it is designed to locally approximate the probabilities assigned to the different classes by the complex model - and this was shown to bear out in experiments over text and images - and it withstood ablation tests. Our experiments, showed that the LIPEX proposal provides more trustworthy and data-efficient explanations compared to multiple other competing methods across various data modalities.

We note that our XAI loss, Equation 6, can be naturally generalized to other probability metrics like the KL divergence. Our studies strongly motivate novel future directions about not only exploring the relative performances between these options but also about obtaining guarantees on the quality of the minima of such novel loss functions.

REFERENCES

- Julius Adebayo, Justin Gilmer, Michael Muelly, Ian Goodfellow, Moritz Hardt, and Been Kim. Sanity checks for saliency maps. *Advances in neural information processing systems*, 31, 2018.
- Sushant Agarwal, Shahin Jabbari, Chirag Agarwal, Sohini Upadhyay, Steven Wu, and Himabindu Lakkaraju. Towards the unification and robustness of perturbation and gradient based explanations. In *International Conference on Machine Learning*, pp. 110–119. PMLR, 2021.
- Raman Arora, Amitabh Basu, Poorya Mianjy, and Mukherjee, Anirbit. Understanding Deep Neural Networks with Rectified Linear Units. In *I.C.L.R.*, 2018.
- Pepa Atanasova, Jakob Grue Simonsen, Christina Lioma, and Isabelle Augenstein. A diagnostic study of explainability techniques for text classification. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pp. 3256–3274, Online, November 2020. Association for Computational Linguistics. doi: 10.18653/v1/2020.emnlp-main.263. URL <https://aclanthology.org/2020.emnlp-main.263>.
- Clément L. Canonne. A short note on learning discrete distributions, 2020.
- Jorge Casillas, Oscar Cordon, Francisco Herrera, and Luis Magdalena. *Accuracy Improvements in Linguistic Fuzzy Modeling*. 01 2003. ISBN 978-3-642-05703-8. doi: 10.1007/978-3-540-37058-1.
- Aditya Chattopadhyay, Kwan Ho Ryan Chan, Benjamin D Haeffele, Donald Geman, and René Vidal. Variational information pursuit for interpretable predictions. *arXiv preprint arXiv:2302.02876*, 2023.
- Joe Collenette, Katie Atkinson, and Trevor Bench-Capon. Explainable ai tools for legal reasoning about cases: A study on the european court of human rights. *Artificial Intelligence*, 317:103861, 2023.
- Jonathan Crabbé and Mihaela van der Schaar. Label-free explainability for unsupervised models. In Kamalika Chaudhuri, Stefanie Jegelka, Le Song, Csaba Szepesvári, Gang Niu, and Sivan Sabato (eds.), *International Conference on Machine Learning, ICML 2022, 17-23 July 2022, Baltimore, Maryland, USA*, volume 162 of *Proceedings of Machine Learning Research*, pp. 4391–4420. PMLR, 2022. URL <https://proceedings.mlr.press/v162/crabbe22a.html>.
- Emmanuel Charleson Dapaah and Jens Grabowski. Debugging machine learning models. 2016.
- Finale Doshi-Velez and Been Kim. Towards a rigorous science of interpretable machine learning, 2017.
- Alexandre Englebort, Olivier Cornu, and Christophe De Vleeschouwer. Poly-CAM: High resolution class activation map for convolutional neural networks, 2022. URL <https://openreview.net/forum?id=qnm-2v-baW>.
- Dumitru Erhan, Yoshua Bengio, Aaron Courville, and Pascal Vincent. Visualizing higher-layer features of a deep network. *University of Montreal*, 1341(3):1, 2009.
- Ruth Fong, Mandela Patrick, and Andrea Vedaldi. Understanding deep networks via extremal perturbations and smooth masks. In *2019 IEEE/CVF International Conference on Computer Vision (ICCV)*, pp. 2950–2958, 2019. doi: 10.1109/ICCV.2019.00304.
- Guang-Hui Fu, Yuan-Jiao Wu, Min-Jie Zong, and Jianxin Pan. Hellinger distance-based stable sparse feature selection for high-dimensional class-imbalanced data. *BMC bioinformatics*, 21(1):1–14, 2020.
- Damien Garreau and Ulrike Luxburg. Explaining the explainer: A first theoretical analysis of lime. In *International conference on artificial intelligence and statistics*, pp. 1287–1296. PMLR, 2020.
- Arushi Gupta and Sanjeev Arora. A simple technique to enable saliency methods to pass the sanity checks, 2020. URL <https://openreview.net/forum?id=BJeGZxrFvS>.

- Arushi Gupta, Nikunj Saunshi, Dingli Yu, Kaifeng Lyu, and Sanjeev Arora. New definitions and evaluations for saliency methods: Staying intrinsic, complete and sound. In S. Koyejo, S. Mohamed, A. Agarwal, D. Belgrave, K. Cho, and A. Oh (eds.), *Advances in Neural Information Processing Systems*, volume 35, pp. 33120–33133. Curran Associates, Inc., 2022. URL https://proceedings.neurips.cc/paper_files/paper/2022/file/d6383e7643415842b48a5077a1b09c98-Paper-Conference.pdf.
- Alon Jacovi, Swabha Swayamdipta, Shauli Ravfogel, Yanai Elazar, Yejin Choi, and Yoav Goldberg. Contrastive explanations for model interpretability, 2021.
- Andrei Kapishnikov, Tolga Bolukbasi, Fernanda Viégas, and Michael Terry. Xrai: Better attributions through regions. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 4948–4957, 2019.
- Andrei Kapishnikov, Subhashini Venugopalan, Besim Avci, Ben Wedin, Michael Terry, and Tolga Bolukbasi. Guided integrated gradients: An adaptive path method for removing noise. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 5050–5058, 2021.
- Zachary Labe and Elizabeth Barnes. Detecting climate signals using explainable ai with single-forcing large ensembles. 04 2021. doi: 10.1002/essoar.10505762.2.
- Himabindu Lakkaraju, Stephen H. Bach, and Jure Leskovec. Interpretable decision sets: A joint framework for description and prediction. In *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, KDD '16, pp. 1675–1684, New York, NY, USA, 2016. Association for Computing Machinery. ISBN 9781450342322. doi: 10.1145/2939672.2939874. URL <https://doi.org/10.1145/2939672.2939874>.
- Benjamin Letham, Cynthia Rudin, Tyler H. McCormick, and David Madigan. Interpretable classifiers using rules and Bayesian analysis: Building a better stroke prediction model. *The Annals of Applied Statistics*, 9(3):1350 – 1371, 2015.
- Ori Linial, Orly Avner, and Dotan Di Castro. Pdexplain: Contextual modeling of pdes in the wild, 2023.
- Scott M. Lundberg and Su-In Lee. A unified approach to interpreting model predictions. In *Proceedings of the 31st International Conference on Neural Information Processing Systems*, NIPS'17, pp. 4768–4777, Red Hook, NY, USA, 2017. Curran Associates Inc.
- Lijun Lyu and Avishek Anand. Listwise explanations for ranking models using multiple explainers. In Jaap Kamps, Lorraine Goeriot, Fabio Crestani, Maria Maistro, Hideo Joho, Brian Davis, Cathal Gurrin, Udo Kruschwitz, and Annalina Caputo (eds.), *Advances in Information Retrieval - 45th European Conference on Information Retrieval, ECIR 2023, Dublin, Ireland, April 2-6, 2023, Proceedings, Part I*, volume 13980 of *Lecture Notes in Computer Science*, pp. 653–668. Springer, 2023.
- Mark S. Neubauer and Avik Roy. Explainable ai for high energy physics, 2022.
- Zhongang Qi, Saeed Khorrarn, and Fuxin Li. Visualizing deep networks by optimizing with integrated gradients, 05 2019.
- Kaivalya Rawal and Himabindu Lakkaraju. Beyond individualized recourse: Interpretable and interactive summaries of actionable recourses. In H. Larochelle, M. Ranzato, R. Hadsell, M.F. Balcan, and H. Lin (eds.), *Advances in Neural Information Processing Systems*, volume 33, pp. 12187–12198. Curran Associates, Inc., 2020.
- Marco Tulio Ribeiro, Sameer Singh, and Carlos Guestrin. "why should i trust you?" explaining the predictions of any classifier. In *Proceedings of the 22nd ACM SIGKDD international conference on knowledge discovery and data mining*, pp. 1135–1144, 2016.
- Marco Tulio Ribeiro, Sameer Singh, and Carlos Guestrin. Anchors: High-precision model-agnostic explanations. AAAI'18/IAAI'18/EAAI'18. AAAI Press, 2018.

- Ramprasaath R Selvaraju, Michael Cogswell, Abhishek Das, Ramakrishna Vedantam, Devi Parikh, and Dhruv Batra. Grad-cam: Visual explanations from deep networks via gradient-based localization. In *Proceedings of the IEEE international conference on computer vision*, pp. 618–626, 2017.
- Avanti Shrikumar, Peyton Greenside, Anna Shcherbina, and Anshul Kundaje. Not just a black box: Learning important features through propagating activation differences. *arXiv preprint arXiv:1605.01713*, 2016.
- Karen Simonyan, Andrea Vedaldi, and Andrew Zisserman. Deep inside convolutional networks: Visualising image classification models and saliency maps, 2014.
- Dylan Slack, Sophie Hilgard, Emily Jia, Sameer Singh, and Himabindu Lakkaraju. Fooling lime and shap: Adversarial attacks on post hoc explanation methods. In *Proceedings of the AAAI/ACM Conference on AI, Ethics, and Society*, AIES '20, pp. 180–186, New York, NY, USA, 2020. Association for Computing Machinery.
- Daniel Smilkov, Nikhil Thorat, Been Kim, Fernanda Viégas, and Martin Wattenberg. Smoothgrad: removing noise by adding noise. *arXiv preprint arXiv:1706.03825*, 2017a.
- Daniel Smilkov, Nikhil Thorat, Been Kim, Fernanda Viégas, and Martin Wattenberg. Smoothgrad: removing noise by adding noise. 06 2017b.
- Kacper Sokol and Peter Flach. Limetree: Consistent and faithful surrogate explanations of multiple classes. *arXiv preprint arXiv:2005.01427*, 2020.
- Jost Tobias Springenberg, Alexey Dosovitskiy, Thomas Brox, and Martin Riedmiller. Striving for simplicity: The all convolutional net. *arXiv preprint arXiv:1412.6806*, 2014.
- Mukund Sundararajan, Ankur Taly, and Qiqi Yan. Axiomatic attribution for deep networks. In *International conference on machine learning*, pp. 3319–3328. PMLR, 2017.
- Berk Ustun, Alexander Spangher, and Yang Liu. Actionable recourse in linear classification. In *Proceedings of the Conference on Fairness, Accountability, and Transparency*, FAT* '19, pp. 10–19, New York, NY, USA, 2019. Association for Computing Machinery. ISBN 9781450361255.
- Sandra Wachter, Brent D. Mittelstadt, and Chris Russell. Counterfactual explanations without opening the black box: Automated decisions and the GDPR. *CoRR*, abs/1711.00399, 2017.
- Ethan Weinberger, Chris Lin, and Su-In Lee. Isolating salient variations of interest in single-cell data with contrastivevi. *Nature Methods*, pp. 1–10, 2023.
- Shawn Xu, Subhashini Venugopalan, and Mukund Sundararajan. Attribution in scale and space. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 9680–9689, 2020.
- Matthew D Zeiler and Rob Fergus. Visualizing and understanding convolutional networks. In *Computer Vision—ECCV 2014: 13th European Conference, Zurich, Switzerland, September 6–12, 2014, Proceedings, Part I 13*, pp. 818–833. Springer, 2014.

A FEATURE SELECTION

Algorithm 1: Forward feature selection

Require: data X , target Y , number of features k

1: $\triangleright X \in \mathbb{R}^{\#\text{Perturbations} \times \#\text{Unique-Words}}$
2: $\triangleright Y \in \mathbb{R}^{\#\text{Perturbations} \times \#\text{Num-Classes}}$

Ensure: set of indices of the selected features Sel_feats

3: $\text{Sel_feats} \leftarrow \{\}$

4: **for** y in Y **do** $\triangleright y \in \mathbb{R}^{\#\text{Perturbations} \times 1}$

5: $f \leftarrow$ initialize selection model of Ridge regression

6: $\text{Current_sel_feats} \leftarrow \{\}$

7: $\text{All_feats} \leftarrow \{1, 2, \dots, \text{len}(X[0])\}$ $\triangleright \text{len}(X[0]) = \text{Number Of Unique Words in } X$

8: **for** $i \leftarrow 1$ to k **do**

9: $\text{best_idx} \leftarrow 0$

10: $\text{best_score} \leftarrow -\infty$

11: **for** $j \in (\text{All_feats} \setminus \text{Sel_feats})$ **do**

12: $f \leftarrow f.\text{fit}(X[:, \text{Sel_feats} \cup \{j\}], y)$

13: $\triangleright f.\text{fit}()$ is used to train f , where the loss function is the linear least squares with l2-norm.

14: $\text{score} \leftarrow$ evaluate f with performance metric of R^2

15: **if** $\text{score} > \text{best_score}$ **then**

16: $\text{best_idx} \leftarrow j$

17: $\text{best_score} \leftarrow \text{score}$

18: **end if**

19: **end for**

20: $\text{Current_sel_feats} \leftarrow \text{Current_sel_feats} \cup \{\text{best_idx}\}$

21: **end for**

22: $\text{Sel_feats} \leftarrow \text{Sel_feats} \cup \text{Current_sel_feats}$

23: **end for**

24: **return** Sel_feats

B LIPEX HYPERPARAMETER SETTINGS

Hyperparameter search was conducted over a small set of randomly selected data of each of the types mentioned below to decide on the following choices.

Learning rate	λ	Batch size
0.01	0.001	128

Table 4: LIPEX Hyperparameter Settings

Note that λ in above refers to the regularizer in the loss in equation 6.

C PSEUDOCODE FOR THE QUANTITATIVE COMPARISON BETWEEN LIPEX AND LIME’S DETECTED IMPORTANT FEATURES (AS GIVEN IN SECTION 4)

Algorithm 2: LIME vs LIPEX w.r.t Angular Spread of the Perturbations About The True Data

Require: k = number of top features to be used for comparing LIME and LIPEX
Require: A set \mathcal{S} of randomly sampled class labelled data at which the comparison is to be done
Require: f^* = the trained predictor that needs explanations.
Require: δ -List of all the angular deviations about the true data at which the LIPEX vs LIME comparison is to be done

- 1: **for** $s \in \mathcal{S}$ **do**
- 2: Compute LIPEX-List- s = top- k features of s w.r.t its predicted class, as detected by the LIPEX matrix using the standard set of Boolean vectors/perturbations w.r.t the all-ones representation of s .
- 3: Compute LIME-List- s = top- k features of s w.r.t its predicted class, as detected by LIME using the standard set of Boolean vectors/perturbations w.r.t the all-ones representation of s - on the same set of features as used in the previous step.
- 4: ▷ **Note that the above two lists of “important” features do not depend on δ ,**
- 5: ▷ **We shall use both as reference lists for the different comparisons to follow.**
- 6: ▷ **The list of features used above will be held fixed in the computations below.**
- 7: **for** $\delta \in \delta$ - List **do**
- 8: Compute δ -LIPEX-List- s = top- k features of s w.r.t its predicted class, as detected by the LIPEX matrix using only those Boolean vectors/perturbations which are within an angle of δ w.r.t the all-ones representation of s .
- 9: Compute δ -LIME-List- s = top- k features of s w.r.t its predicted class, as detected by LIME using only those Boolean vectors/perturbations which are within an angle of δ w.r.t the all-ones representation of s
- 10: Compute the Jaccard Index, $J_{s,\delta\text{-LIPEX-vs-LIME}} := \frac{|\delta\text{-LIPEX-List-}s \cap \text{LIME-List-}s|}{|\delta\text{-LIPEX-List-}s \cup \text{LIME-List-}s|}$
- 11: Compute the Jaccard Index, $J_{s,\delta,\text{LIME}} := \frac{|\delta\text{-LIME-List-}s \cap \text{LIME-List-}s|}{|\delta\text{-LIME-List-}s \cup \text{LIME-List-}s|}$
- 12: Compute the Jaccard Index, $J_{s,\delta,\text{LIPEX}} := \frac{|\delta\text{-LIPEX-List-}s \cap \text{LIPEX-List-}s|}{|\delta\text{-LIPEX-List-}s \cup \text{LIPEX-List-}s|}$
- 13: **end for**
- 14: **end for**
- 15: **end for**
- 16: Plot $\left(\frac{1}{|\mathcal{S}|} \cdot \sum_{s \in \mathcal{S}} J_{s,\delta\text{-LIPEX-vs-LIME}} \right)$ vs δ
- 17: Plot $\left(\frac{1}{|\mathcal{S}|} \cdot \sum_{s \in \mathcal{S}} J_{s,\delta,\text{LIME}} \right)$ vs δ
- 18: Plot $\left(\frac{1}{|\mathcal{S}|} \cdot \sum_{s \in \mathcal{S}} J_{s,\delta,\text{LIPEX}} \right)$ vs δ

D ADDITIONAL EXPERIMENTS

D.1 δ EFFECT ON JACCARD EXPERIMENT

For Text data				
δ (radians)	$\frac{\pi}{16}$	$\frac{\pi}{8}$	$\frac{\pi}{4}$	$\frac{\pi}{2}$
number of perturbation points	138	659	2383	5000
For Image data				
δ (radians)	$\frac{7\pi}{30}$	$\frac{8\pi}{30}$	$\frac{9\pi}{30}$	$\frac{10\pi}{30}$
number of perturbation points	228	774	994	1000

Table 5: The effect of δ on the number of perturbation points, result averaged on 100 input instances.

Note that when δ decreases, while the amount of allowed perturbations falls, the similarity measure π in equation 6 increases.

D.2 ADDITIONAL DATA FOR THE DEMONSTRATION IN FIGURE 1

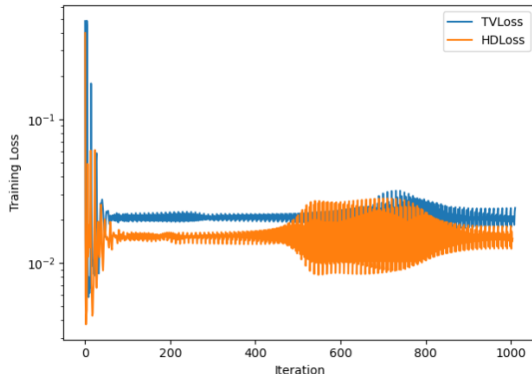


Figure 5: Comparison of the progress of training on the Hellinger distance based LIPEX loss, as given in equation 6, (its training curve being labelled as “HDLoss” above) and a natural Total Variation distance analogue of it (its training curve labelled as “TVLoss” above), for the text data in Figure 1

D.3 LIPEX ON IMAGE

Each class/row of our explanation matrix would contain a weight corresponding to the importance of a common set of features/super-pixels for that class. The figure below shows the part of the matrix corresponding to the top 3 classes detected for this image i.e. “Burmese_mountain_dog”, “Entlebucher” and “Appenzeller” and the top-4 features deemed to be important for the predicted class i.e. “Burmese_mountain_dog”. Thus we see how LIPEX successfully “localized” the dog as being determinant to the predictions rather than the cat which is also prominent in this picture.

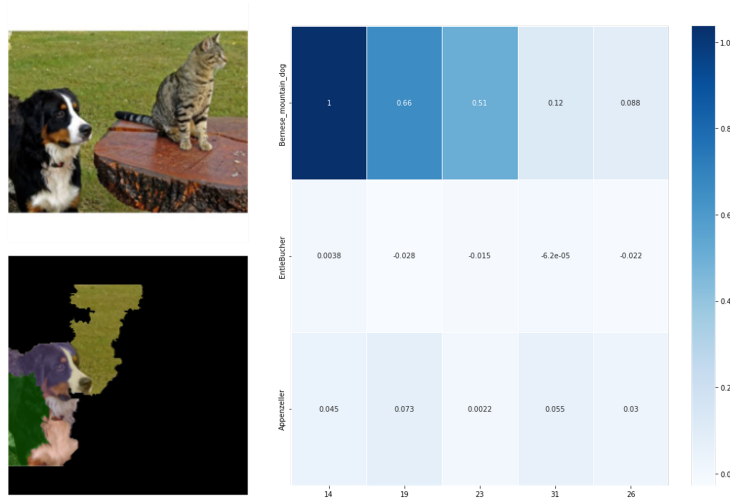


Figure 6: Example of (the most important part) of the matrix returned by the LIPEX method on an image. See Figure 7 how the top 5 segments detected for the image patch together. The corresponding LIME answer is visualized in Figure 8 - and we can see how it prioritized image segments unrelated to the dog.

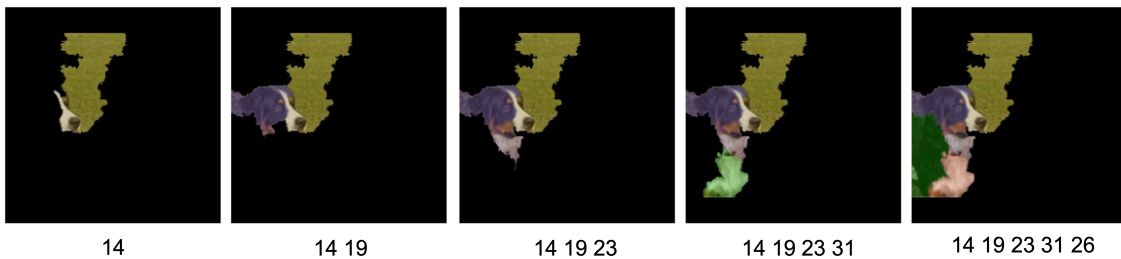


Figure 7: The top 5 image segments deemed to be important by LIPEX for Inception-V3 to classify the image in Figure 6 as a "Burmese_mountain_dog"



Figure 8: The weight vector over the features as returned by LIME

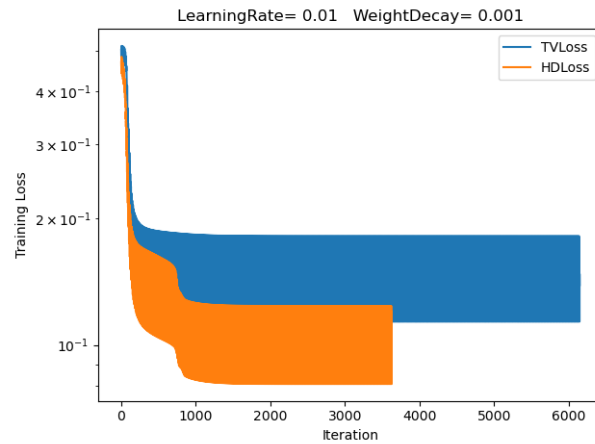


Figure 9: The LIPEX result for the image data in Figure 7 corresponds to the endpoint of minimizing the loss given in equation 6. (its training curve being labelled as “HDLoss” above). The training curve labelled as “TVLoss” corresponds to training on an analogous loss as in equation 6 but with the metric used in probability space being Total Variation.

D.4 MORE DETAILS ABOUT THE ABLATION STUDIES IN SECTION 4

Method	Top3	Top4	Top5
Vgg16			
LIPEX	0.763 (± 0.026)	0.82 (± 0.014)	0.867 (± 0.017)
LIME	0.74(± 0.014)	0.78(± 0)	0.793(± 0.009)
XRAI	0.713(± 0.025)	0.77(± 0.008)	0.793(± 0.026)
GradCAM	0.69(± 0.016)	0.763(± 0.026)	0.79(± 0.043)
GuidedIG	0.717(± 0.017)	0.75(± 0.016)	0.793(± 0.009)
BlurIG	0.713(± 0.009)	0.787(± 0.017)	0.807(± 0.017)
Vanilla_Grad	0.68(± 0.022)	0.753(± 0.005)	0.807(± 0.012)
SmoothGrad	0.747(± 0.019)	0.817(± 0.005)	0.843(± 0.005)
Integrated_Grad	0.703(± 0.021)	0.747(± 0.029)	0.773(± 0.017)
InceptionV3			
LIPEX	0.673(± 0.005)	0.753 (± 0.005)	0.773 (± 0.017)
LIME	0.63(± 0.014)	0.713(± 0.046)	0.767(± 0.034)
XRAI	0.693 (± 0.017)	0.7(± 0.029)	0.74(± 0.062)
GradCAM	0.653(± 0.005)	0.697(± 0.017)	0.72(± 0.036)
GuidedIG	0.663(± 0.046)	0.67(± 0.051)	0.713(± 0.04)
BlurIG	0.647(± 0.041)	0.703(± 0.034)	0.717(± 0.029)
Vanilla_Grad	0.657(± 0.012)	0.653(± 0.026)	0.72(± 0.051)
SmoothGrad	0.65(± 0.024)	0.683(± 0.049)	0.74(± 0.037)
Integrated_Grad	0.637(± 0.019)	0.707(± 0.025)	0.733(± 0.04)

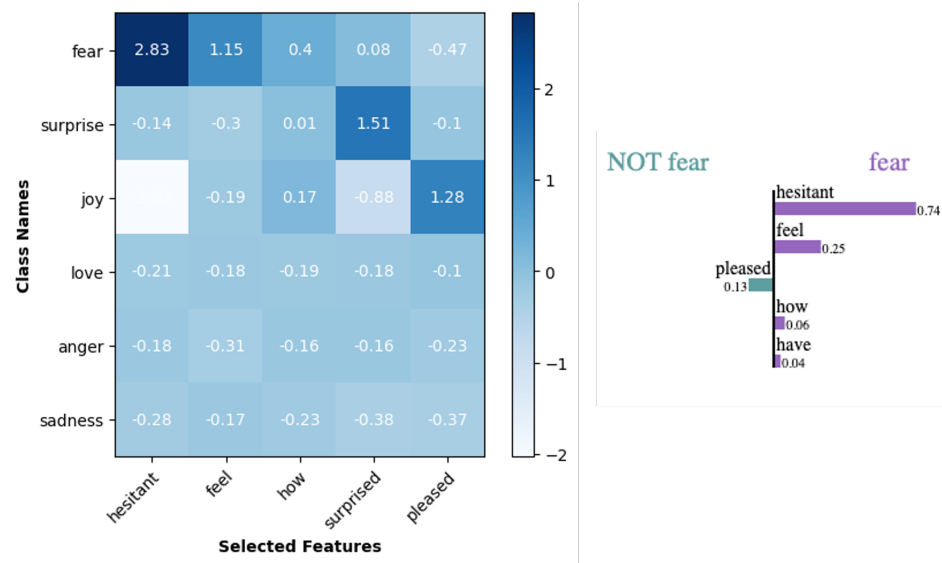
Table 6: Ablation study of Vgg16 and InceptionV3, results averaged on three times run, each run over randomly chosen 100 images from the ImageNet validation dataset. Here, features refer to image segments (via Segment Anything) weighted by LIME/LIPEX or pixel importances determined by saliency methods.

Method	Top1	Top2	Top3	Top4	Top5
BERT on 20Newsgroups					
LIPEx	0.781 (±0.047)	0.857 (±0.036)	0.897 (±0.016)	0.909 (±0.026)	0.908(±0.032)
LIME	0.777(±0.027)	0.841(±0.031)	0.856(±0.045)	0.881(±0.011)	0.912(±0.021)
GuidedBack	0.387(±0.049)	0.477(±0.082)	0.517(±0.074)	0.517(±0.054)	0.553(±0.042)
Saliency	0.387(±0.049)	0.477(±0.082)	0.517(±0.074)	0.517(±0.054)	0.553(±0.042)
Input_G	0.38(±0.054)	0.48(±0.079)	0.52(±0.071)	0.523(±0.05)	0.57(±0.029)
Deeplift	0.38(±0.054)	0.48(±0.079)	0.52(±0.071)	0.523(±0.05)	0.57(±0.029)
Occlusion	0.45(±0.045)	0.543(±0.066)	0.59(±0.082)	0.627(±0.065)	0.653(±0.063)
BERT on Emotion					
LIPEx	0.657 (±0.021)	0.74 (±0.037)	0.73 (±0.028)	0.73 (±0.029)	0.793 (±0.024)
LIME	0.653(±0.017)	0.697(±0.041)	0.647(±0.037)	0.64(±0.045)	0.65(±0.008)
GuidedBack	0.597(±0.029)	0.61(±0.029)	0.637(±0.009)	0.623(±0.009)	0.63(±0.008)
Saliency	0.597(±0.029)	0.61(±0.029)	0.637(±0.009)	0.623(±0.009)	0.63(±0.008)
Input_G	0.6(±0.033)	0.62(±0.029)	0.637(±0.009)	0.633(±0.009)	0.637(±0.005)
Deeplift	0.6(±0.033)	0.63(±0.028)	0.653(±0.012)	0.643(±0.012)	0.64(±0.008)
Occlusion	0.65(±0.016)	0.66(±0.022)	0.697(±0.012)	0.697(±0.005)	0.693(±0.017)

Table 7: Ablation study on Text dataset by removing TopK features detected by explainable methods and doing re-prediction, each experiment was independently repeated three times on randomly chosen 100 text instances.

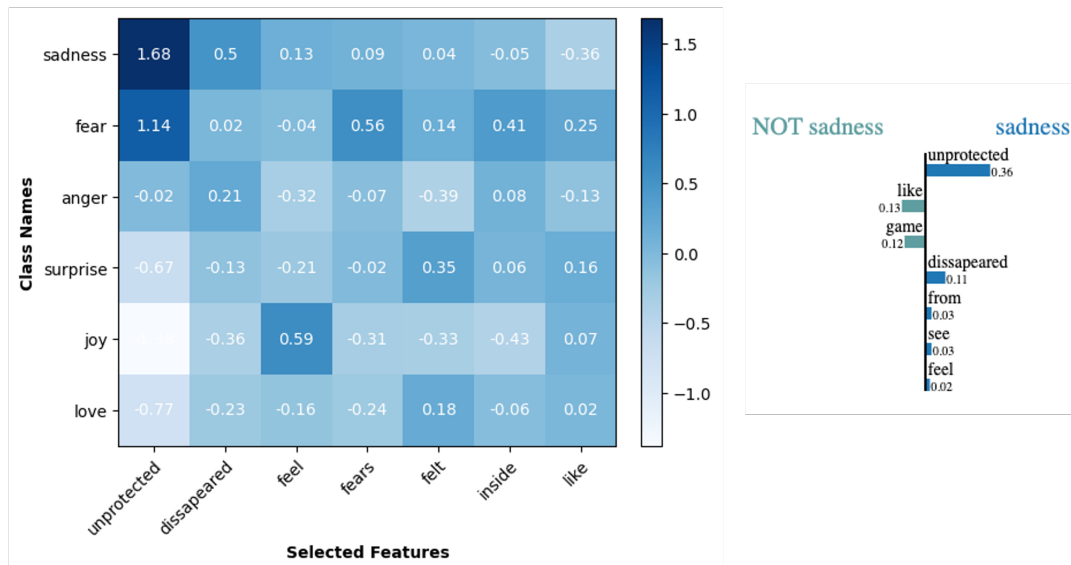
D.5 COMPARISON BETWEEN THE EXPLANATIONS FOUND BY LIPEX AND LIME

Figure 10, 11 vividly demonstrate the fine-grained explanation that is obtained by the LIPEX matrix as opposed to the LIME’s explanation (right bar) on the same number of feature sets. The input instances in Figure 10 and 11 are randomly chosen from the Emotion dataset. The explanatory matrix generated by LIPEX makes it easy to see the relationship between the same feature and different categories.



Text: *i feel like in the last year especially i ve gone from a girl to a woman and despite how hesitant i have always been about getting older next year i will be twenty four i am surprised at how pleased i am to have done so.*

Figure 10:



Text: *i feel inside this life is like a game sometimes then you came around me the walls just dissappeared nothing to surround me keep me from my fears im unprotected see how ive opened up youve made me trust coz ive never felt like this before im naked around you does it show.*

Figure 11:

D.6 COMPARISON BETWEEN THE EXPLANATIONS FOUND BY LIPEX AND LIME FOR TEXT DATA WHERE PREDICTED CLASS AND THE TRUE CLASS ARE DIFFERENT

Figure 12, 13, 14, 15 demonstrate the fine-grained explanation that is obtained by the LIPEX matrix as opposed to the LIME’s explanation on the same feature set for the predicted class - which is different than the true class for these instances. The input instances in Figure 12, 13, 14 and 15 are chosen from 20Newsgroups.

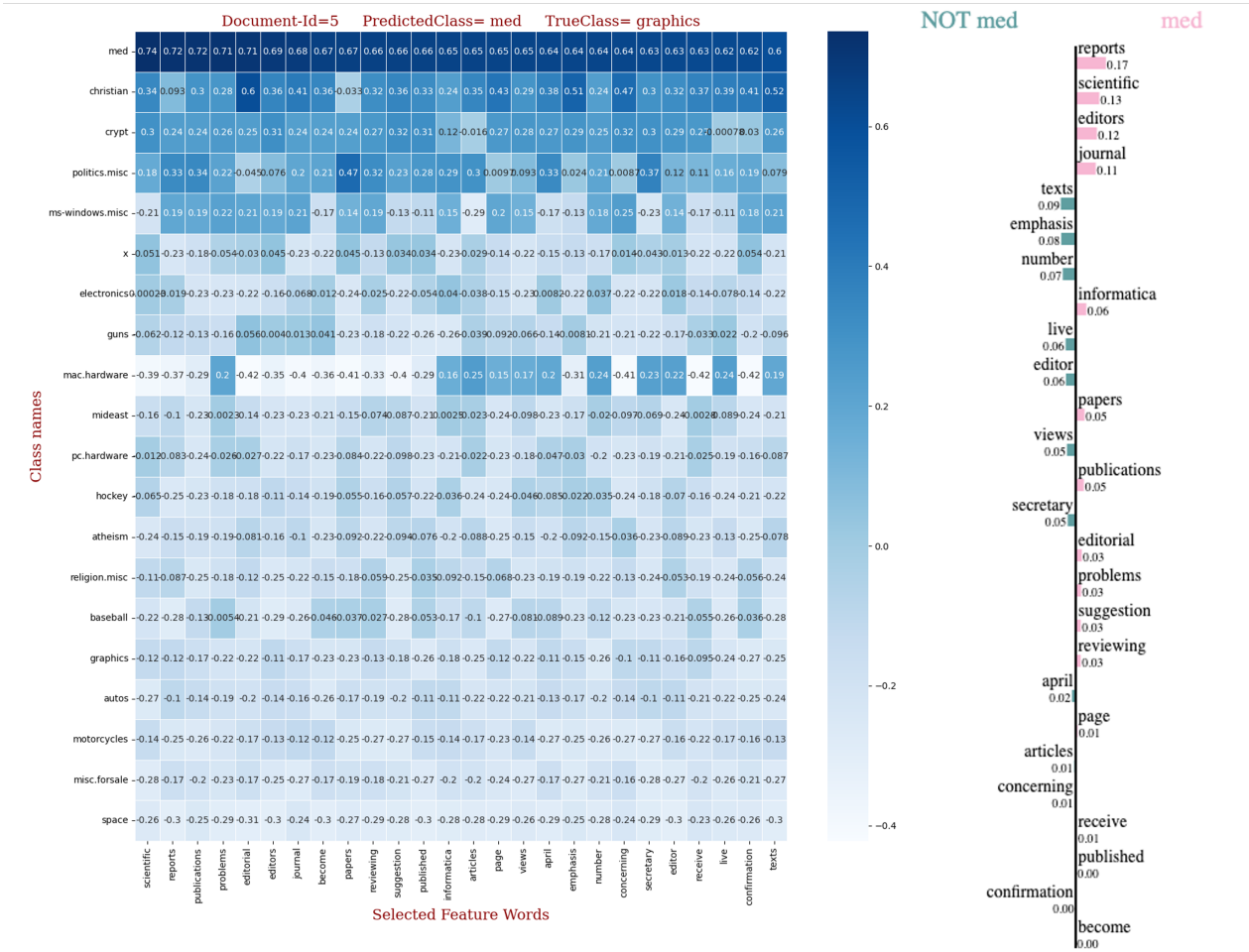


Figure 12:

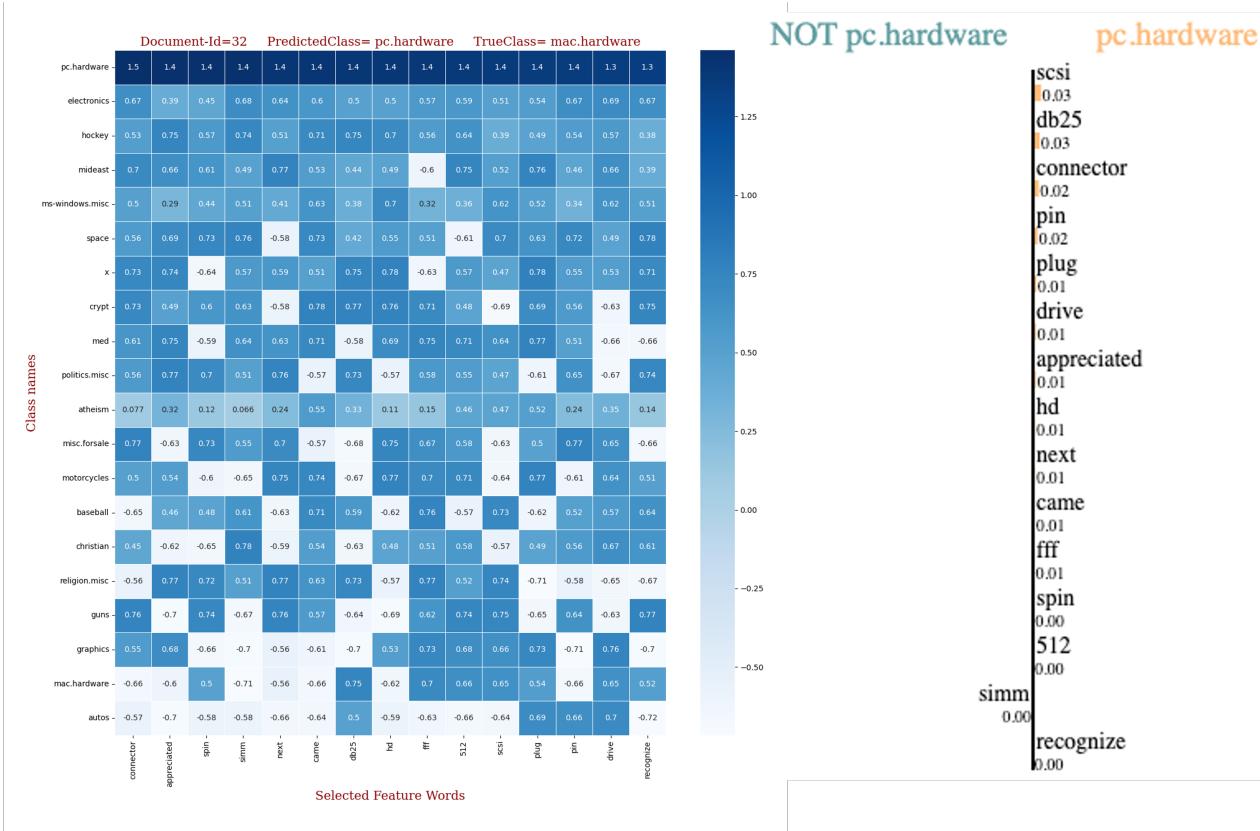


Figure 13:

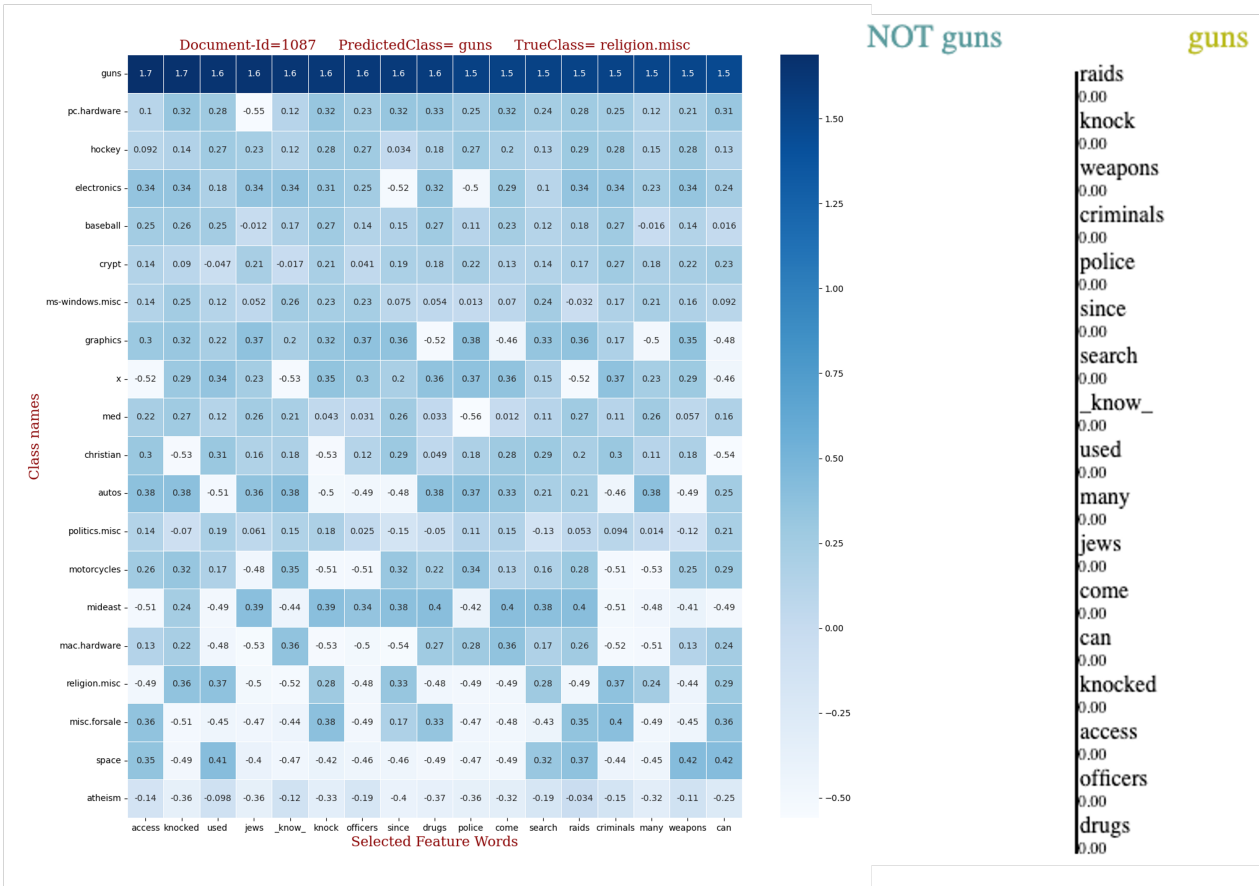


Figure 14:



Figure 15: