

Classifier Data Quality – A Geometric Complexity Based Method for Automated Baseline And Insights Generation

George Kour^{1,2}, Marcel Zalmanovici¹, Orna Raz¹, Samuel Ackerman¹, Ateret Anaby-Tavor¹

¹IBM Research, Haifa

²Sagol Department of Neurobiology, University of Haifa

Abstract

Testing Machine Learning (ML) models and AI-Infused Applications (AIAs), or systems that contain ML models, is highly challenging. In addition to the challenges of testing classical software, it is acceptable and expected that statistical ML models sometimes output incorrect results. A major challenge is to determine when the level of incorrectness, e.g., model accuracy or F1 score for classifiers, is acceptable and when it is not. In addition to business requirements that should provide a threshold, it is a best practice to require any proposed ML solution to out-perform simple baseline models, such as a decision tree.

We have developed complexity measures, which quantify how difficult given observations are to assign to their true class label; these measures can then be used to automatically determine a baseline performance threshold. These measures are superior to the best practice baseline in that, for a linear computation cost, they also quantify each observation's classification complexity in an explainable form, regardless of the classifier model used. Our experiments with both numeric synthetic data and real natural language chatbot data demonstrate that the complexity measures effectively highlight data regions and observations that are likely to be misclassified.

1 Introduction

Testing AIAs is highly challenging. The characteristics of the data used for training an ML model are key to the quality of the resulting ML model. The geometric foundation underlying many classifiers is to assign an observation to the closest class. However, different classifiers may use different geometrical properties to capture the class geometry as well as may employ different distance functions to capture the closeness between an observation and a class geometry. Here, we present a complexity heuristic based on this intuition that attempts to quantify the difficulty of assigning a given observation to its true class label by calculating the relative closeness of the observation to its true class vs. its closeness to other classes.

The complexity measure provides insights for improving the system design, for improving labeling, or for directing the gathering or generation of more data. It can further be used to set an expected performance level, which we call a

'baseline', for a simple distance-based discriminator on the entire dataset. Section 2 provides more details. Our complexity measures have the following properties:

- They are computationally efficient and feasible to calculate quickly.
- They provide information both about single observations and about sets of observations, such as those that have the same label or class.
- They are general and are applicable to a wide variety of discrimination algorithms as they share similar underlying geometrical properties, stemming from Gaussian Discriminant Analysis, GDA, (Ghojogh and Crowley 2019).
- They are explainable in that the geometry provides a clear reason for indicating an observation or a set of observations as complex.

Our method overcomes some of the weaknesses of the best practice baseline: it uses all the data and does not necessitate dividing the data into train and validation sets — this is especially important when the data is small; it highlights difficult to classify observations at the same low calculation cost — compared to manual debugging and very costly methods for scoring individual observations such as leave-one-out and Shapley values; it provides a graphical explanation of the scores (see examples in Section 3).

We assess the properties and effectiveness of our complexity measures over numeric synthetic data. Then we demonstrate the effectiveness of our measures over real natural language data. We focus on the especially challenging domain of chatbots. Chatbots are a prominent example of an AIA. Chatbots allow users to interact with the business through a natural language interface and are becoming a key channel for customer engagement. For many customers, the chatbot provides their first interaction with the business (e.g. for support purposes). Therefore, it is important for a chatbot to be of good quality from day one. Chatbot technology usually comprises two basic components:

1. A machine learning (ML) natural language processing (NLP) based intent classifier that can process what the user is saying, and
2. A conversation flow orchestrator that incorporates domain knowledge and is driven by the business actions and

potentially content extracted from past human-to-human dialogs and company documents.

We utilize our complexity measures to assess the quality of the intent classifier data and to set a baseline for any intent classifier trained with that data.

As part of our methodology we automatically find the interesting measures and ranges of complexity. These are the ranges that are more likely to result in misclassification. For that purpose we utilize our IBM FreaAI technology (Barash et al. 2019; Ackerman, Raz, and Zalmanovici 2020). Automatically highlighting these ranges assists in fault identification as well as in selecting additional data for training or for testing. This data may be automatically generated, e.g., through the LAMBADA technology (Anaby-Tavor et al. 2020).

The rest of the paper is organized as follows: Section 2 introduces our complexity measures and their usage in automatically generating a baseline and indicating complex observations that are likely to challenge many discrimination algorithms. Section 3 demonstrates how we utilize these measures for setting a baseline by identifying those observations that are likely to be misclassified. Section 4 summarizes related work and Section 5 concludes and discusses limitations of our approach.

2 Methodology

In this section we describe our framework for estimating the complexity—in terms of difficulty of assignment to its true label—of a sample in a labeled data set or corpus, in the case of text data. We also provide a methodology for using the resulting complexity measures when automatically setting a baseline for many common discrimination algorithms as well as providing additional insights about observations and their complexity. Complex to classify observation are highlighted in an explainable way. These may be observations that are part of the data used to compute the complexity measures, such as those in the train set, or observations yet unseen, such as those in the test set or those observations that came at runtime after deployment (e.g., user utterances after the chatbot is deployed).

Given a labeled corpus, we first employ a text embedding model e_θ on each of the samples to obtain a dataset $\mathcal{D} = \{(x_i, y_i)\}_{i=1}^N$ where $x_i = e_\theta(s_i)$ and $x_i \in \mathcal{X} \subseteq \mathbb{R}^d$ is a d -dimensional dense vector representation of the textual sample s_i , and $y_i \in \mathcal{Y}$ is its corresponding label. In this work, we employ the SENTENCE-TRANSFORMERS (Reimers and Gurevych 2019) python package and the attention-based pre-trained language model called PARAPHRASE-MINI-LM-L6-v2 (Reimers and Gurevych 2019) which maps textual samples to \mathbb{R}^{384} .

Our goal is to define a simple yet effective geometric-based measure $h : \mathcal{X} \times \mathcal{Y} \rightarrow \mathbb{R}_{\geq 0}$ which evaluates the incompatibility of a given sample (x, y) to the underlying distribution and thus estimating how difficult it would be for a discriminator to correctly predict its label. In essence, the idea behind our complexity measure is to calculate the relative geometrical closeness of the sample to its own class versus its closeness to other classes in the dataset. The relevant

geometric properties considered in the calculation are determined by the employed distance measure $\delta(x, B)$ which should capture the perceptual distance between the sample x and the geometry of the class samples in B . Thus, we define the complexity heuristic as follows:

$$\begin{aligned} h(x, y) &= -\log \left[\frac{e^{-\delta(x, \mathcal{C}(y))}}{\sum_{c \in \mathcal{Y}} e^{-\delta(x, \mathcal{C}(c))}} \right] \\ &= \delta(x, y) + \log \left[\sum_{c \in \mathcal{Y}} e^{-\delta(x, \mathcal{C}(c))} \right] \end{aligned} \quad (1)$$

where $\mathcal{C}(c)$ represents the set of the samples in \mathcal{D} that belong to class c (namely, $\mathcal{C}(c) = \{x | (x, c) \in \mathcal{D}\}$).

We compare multiple distance functions which capture different geometric properties: Euclidean (E), Cosine similarity (S), and Mahalanobis (M) (McLachlan 1999), defined respectively in Equations 2, 3 and 4 below. Note that while the Euclidean and the cosine metrics consider only the classes’ centroids, the Mahalanobis distance captures the dispersion of the classes along each of the embedding dimensions, and calculates the “effective” distance of the sample from the class geometry.

$$\delta_E(a, B) = \|a - \hat{\mu}_B\|_2 \quad (2)$$

$$\delta_S(a, B) = \frac{a \cdot \hat{\mu}_B}{\|a\| \|\hat{\mu}_B\|} \quad (3)$$

$$\delta_M(a, B) = \sqrt{(a - \hat{\mu}_B)^T \hat{\Sigma}_B^{-1} (a - \hat{\mu}_B)} \quad (4)$$

where B is a class of samples, and $\hat{\mu}_B$ and $\hat{\Sigma}_B$ are the set centroid and empirical covariance matrix, calculated as follows:

$$\hat{\mu}_B = \frac{1}{|B|} \sum_{b \in B} b \quad (5)$$

$$\hat{\Sigma}_B = \frac{1}{|B|} \sum_{b \in B} (b - \hat{\mu}_B)(b - \hat{\mu}_B)^T \quad (6)$$

Using Mahalanobis distance in the complexity measure is equivalent to calculating negative log-likelihood loss of a generative classification model which fits a class-conditioned Gaussian distribution to every class of the data. In (Lee et al. 2018) a similar method is used to calculate confidence scores for out-of-distribution (OOD) detection in test samples, however, in their case the confidence function $M(x)$ is independent of the given sample label and thus for a given feature vector x , they defined the confidence score as follows:

$$M(x) = \max_c \{-\delta_M(x, \mathcal{C}(c))\} \quad (7)$$

Despite being an asymptotically unbiased estimator of the covariance matrix, the Maximum Likelihood Estimator (Equation 6) is imperfect and the precision matrix obtained from its inversion may be inaccurate and sometimes even unobtainable due to numerical reasons. Estimating the covariance matrix is even harder in cases where the number of samples is smaller than the number of features (i.e., $n < p$).

Thus, in relatively small datasets, the empirical covariance matrix is commonly underestimated. To mitigate the estimation error of the covariance matrix in high-dimensional feature spaces shrinkage-based covariance estimation methods, such as the Ledoit-Wolf shrinkage approach (Ledoit and Wolf 2004), are commonly employed.

We found that using the estimated Pearson correlation coefficients matrix (Galton 1877) instead of the covariance matrix in the Mahalanobis distance for calculating the complexity yields better results, in the sense that it better correlates with the error concentration of a prototypical classifier (more information is provided in Section 3). In the following, using the correlation matrix is denoted by $\delta_{\hat{M}}$. We hypothesise that despite the lose of information about the manifold radii, using correlation coefficients matrix is the better option in this case as it is less affected by the number of samples.

The complexity measure may be employed on both the train and test sets. While its goal in both cases is to identify difficult to classify samples, its practical interpretation can be different depending on the usage setup. For instance, when used on the train set it may be used to remove training samples which may be mis-labeled or spot samples with high information gain; however, when used on the test set, high complexity samples can be viewed as ambiguous or out-of-domain. Note that when the complexity measure is used to rank the train samples, their classification difficulty may be underestimated because the estimated geometry of the classes is affected by the samples to be ranked. However, when having enough samples, the marginal effect of each single sample on the estimated geometry is small.

2.1 Setting a baseline

In our framework, given a dataset \mathcal{D} and a distance function $\delta(x, B)$, where x is a sample in the d -dimensional space and B is a set of samples, the baseline classifier is defined as follows:

$$f(x) = \operatorname{argmin}_{c \in \mathcal{Y}} \delta(x, C(c)) \quad (8)$$

Note that the complexity measure is defined as the negative log-likelihood of $f(x)$. Using δ_E or δ_S distance functions, the classifier becomes a Nearest-Centroid classifier in terms of amplitude or angle, respectively. Observations that are more difficult to classify by class-centroid distance will receive higher complexity scores under our heuristic.

However, when the Mahalanobis distance is used, i.e., this classifier becomes a Gaussian Discriminant Analysis (GDA) classifier, and when the Pearson correlation coefficient is used instead of the covariance matrix then the classifier is reduced to Quadratic Discriminant Analysis classifier (QDA). These classifiers are attractive because they have been proven to work, have no hyperparameters to tune, and have closed-form solutions that can be easily computed.

3 Experiments and results

Evaluating the complexity measure involves assessing its correlation with both the errors and confidence of a somehow prototypical classifier. Namely, high complexity values

are expected for samples on which the classifier tends to err or for samples that may be classified correctly by the classifier but with low confidence. Moreover, when a trained classifier errs with high confidence, the complexity measure may correctly indicate high complexity.

Our experimental results indicate that the complexity measures are similar to or even out-perform a trained classifier confidence measure. This is especially encouraging as it means that the complexity measures can reasonably estimate the difficulty a classifier would have to correctly classify samples, even prior to any training.

We experiment with both synthetic numeric data as described in Section 3.1 and with real natural language chatbot data as described in Section 3.2. We utilize FreaAI for automatically finding the most correlating features and ranges with high misclassification concentration, which we refer to as *Slices*. For each slice computed by FreaAI, a relative rank is calculated, which can be intuitively interpreted as follows: the higher the ranking, the higher the error concentration of the classifier on that slice compared to overall performance. In addition to the performance of the classifier in the slice, the ranking also considers the number of samples in the dataset that fall into the slice.

3.1 Synthetic data

Note that for different needs, the complexity measure can be employed on either the train samples or on other samples not used for training. In the train use-case case, the geometries and classifier are learned approximated using the train set, and the complexity is calculated on the same train samples. However, in the test use-case, the training and the approximation of the geometries is performed using the train set, but the complexity is calculated on the new samples.

To demonstrate the complexity measure potency in both the train and test use-cases, our experiments include calculating the complexity measure on the train set of synthetic 2-dimensional data, and on the test set of a real-life sentence utterance classification (SUC) datasets.

In the first set of experiments, we generated synthetic samples for three classes c_0, c_1, c_2 by drawing random samples from three bivariate normal distributions. The distributions are specified by their mean and covariance matrix. Nevertheless, the classes geometries partially overlap to simulate ambiguous samples. Second, the complexity measure is calculated for each sample and SVM classifier (H. et al. 1996) is trained to predict the sample classes. Then, the classifier is used to predict the samples labels and the prediction is analyzed using FreaAI technology to find and rank slices with high error concentrations. Over-fitting of the classifier is avoided by applying a relatively simple model, and thus, we expect most errors in predicting the labels to be in the overlap. The advantage of using synthetic datasets is that it makes the analysis and visual validation easier. The main purpose is to make sure that the complexity highlights areas where it is more likely to encounter errors and that the FreaAI technology identifies the expected complexity ranges as problematic. The data records sent to FreaAI contain for each sample its coordinates x_1, x_2 , the true label y , the label predicted by the trained classifier \hat{y} and its prediction confi-

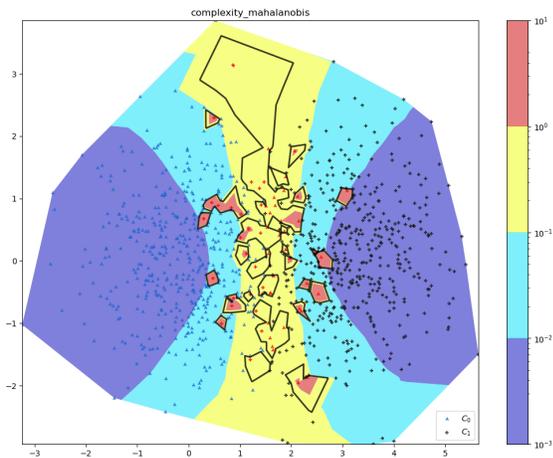


Figure 1: Synthetic data drawn from two Gaussians of equal size; background the complexity calculated via the Mahalanobis method. In all Figures the background shows the complexity areas that were interpolated by the samples in the dataset. The color legend goes from most complex at the top to least complex at the bottom of the legend. The samples and their class label are represented by colored points. The areas bounded by the black contour represent slices found by FreaAI and contain high concentration of errors.

dence as well as the different complexity measures proposed in Section 2.

In the first and simplest experiment we use only two classes of the same size that overlap as seen in Figure 1. In this case we expect most errors to be located in the overlap as well as for those samples to be scored with higher complex. The generated data has 500 samples in each class and 63 errors in total (accuracy 0.94). The FreaAI slices on the Mahalanobis and Euclidean complexity confirm that this is where the errors are found as can be seen in Figure 1.

The background of all figures with a complexity measure is similar to Figure 1 and shows the complexity areas that were interpolated by the samples in the dataset. The color legend goes from most complex at the top to least complex at the bottom of the legend. The samples and their class label are represented by colored points. The areas bounded by the black contour represent slices found by FreaAI and contain high concentration of errors.

All experiments similarly capture the higher complexity ranges that correlate best with misclassifications. Figure 1 and its corresponding Table 1 mark, for example, a range where 65 records appear which contains all 63 errors in this data. The cosine complexity measure marked a much larger range and still didn't cover most errors. The model confidence in this case shows that all errors are when the confidence is less than 0.96 which does not help much. We show also a few interesting ranges found by FreaAI on x_1 , x_2 or both where relatively many errors are found. Not surprisingly these areas are all in the overlap between the classes.

For the second experiment we turned one of the classes into an ellipse as depicted in Figure 2. This is to verify that

Feature(s)	Slice	acc_{SVC}	Size	Rank
compl_mah	0.46 ~ 5.40	0.05	65	0.99
compl_euc	0.46 ~ 5.40	0.05	65	9.98
x_1	0.30 ~ 2.57	0.81	330	0.58
confidence	0.50 ~ 0.96	0.80	309	0.56
compl_cos	0.34 ~ 1.61	0.84	168	0.55
x_1, x_2	1.52 ~ 1.69	0.26	19	0.54
	-1.43 ~ 1.80			
x_1, x_2	1.20 ~ 1.44	0.27	11	0.35
	-0.54 ~ 1.26			
compl_cos	2.38 ~ 2.41	0.58	12	0.15
x_2	0.85 ~ 1.89	0.71	17	0.13
x_2	-0.54 ~ -0.51	0.54	11	0.10

Table 1: Slices for two Gaussians of equal size. These slices are depicted in Figure 1 as areas bounded by black contour.

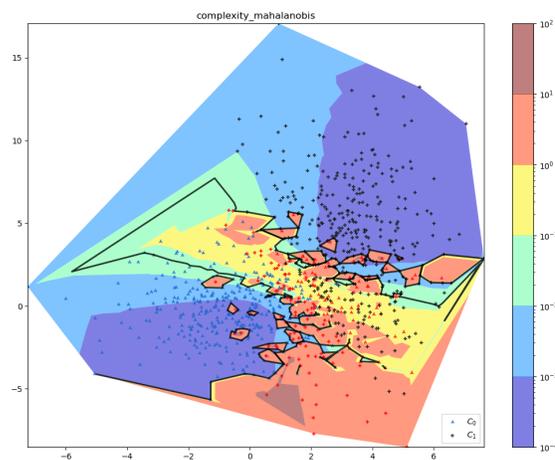


Figure 2: Similar to Figure 1 however here the shape of one class is still a circle while the other is an ellipse.

the complexity works as expected when the distance on each dimension varies. As can be seen, the higher complexity is marked in the main intersection of the Gaussians and also in the bottom-right and top-left sides where the narrower ellipse points are drawn. This data has 181 errors (accuracy 0.82). For example the Mahalanobis complexity range identified by FreaAI contains 141 of those errors, versus 371 found in the Euclidean complexity and 76 in the "best" confidence range.

For the third experiment we generated three same-sized Gaussians that overlap in roughly the same area. This data contains 160 mistakes (accuracy 0.89). Figure 3 shows the Euclidean complexity as the background for the data. This is almost identical to the heat map generated by the Mahalanobis complexity and quite similar to the model confidence heat map shown in Figure 4.

In the fourth experiment we used three classes: two elliptic Gaussians and one round. To make it more interesting we also made sure that x_1 overlapped with x_2 in one area and with x_3 in another. The result can be seen in Figures 5–4.

Feature(s)	Slice	acc_{SVC}	Size	Rank
compl_mah	0.04 ~ 17.87	0.45	320	0.51
compl_euc	0.03 ~ 18.63	0.49	351	0.48
x_1, x_2	0.34 ~ 5.30	0.36	47	0.47
	-8.53 ~ -2.88			
x_1, x_2	0.82 ~ 4.27	0.63	323	0.45
	-2.81 ~ 3.14			
compl_cos	1.32 ~ 2.83	0.55	182	0.42
confidence	0.50 ~ 0.63	0.50	151	0.40
x_1	1.80 ~ 2.60	0.61	137	0.38
x_2	8.53 ~ 3.02	0.55	64	0.32

Table 2: Slices for two Gaussians of unequal size as shown in Figure 2.

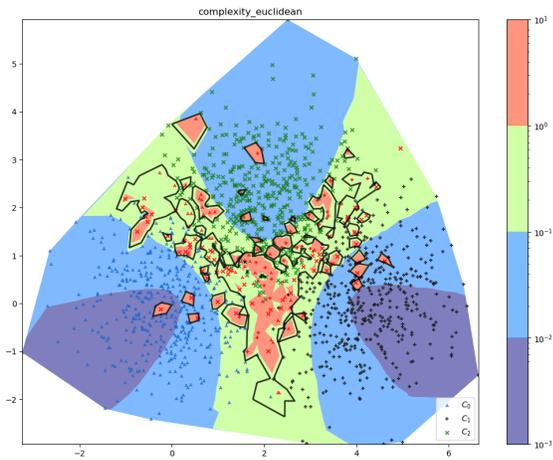


Figure 3: Synthetic data with three Gaussians overlapping in the same area; background shows the Euclidean complexity.

Feature(s)	Slice	acc_{SVC}	Size	Rank
compl_euc	0.54 ~ 8.93	0.16	187	0.93
compl_mah	0.54 ~ 8.98	0.16	187	0.92
compl_cos	0.93 ~ 3.58	0.76	485	0.49
confidence	0.37 ~ 0.86	0.66	426	0.48
x_1, x_2	-0.98 ~ 2.91	0.73	185	0.45
	0.62 ~ 1.42			
x_2	0.87 ~ 1.36	0.74	173	0.44
x_1, x_2	2.98 ~ 4.22	0.65	95	0.40
	0.64 ~ 2.27			
x_1	1.45 ~ 1.58	0.76	46	0.32

Table 3: Slices for three Gaussians of equal size. (Some) of the results of the data shown in Figure 3 and 4.

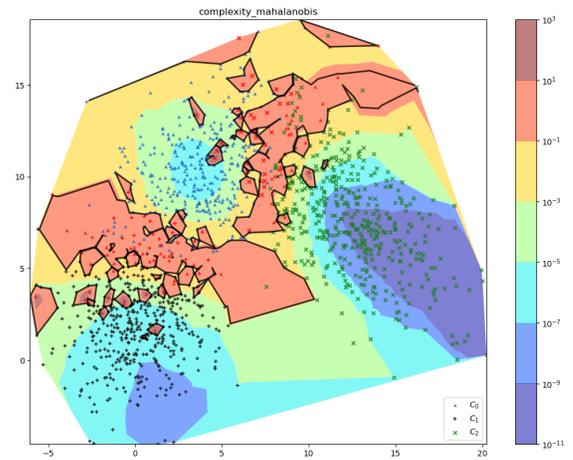


Figure 5: Three Gaussians with two distinct overlapping areas; background is Mahalanobis complexity.

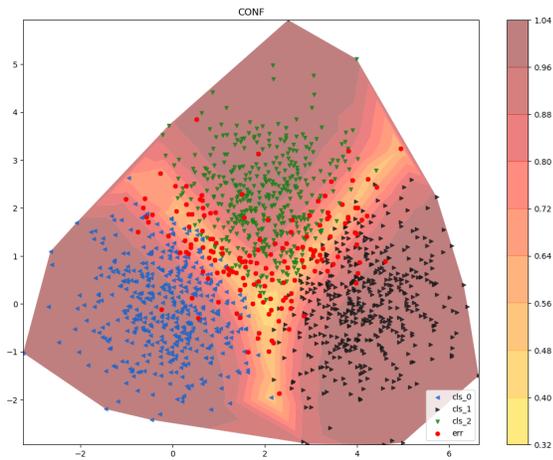


Figure 4: The same data as in Figure 3 with the background showing the model confidence.

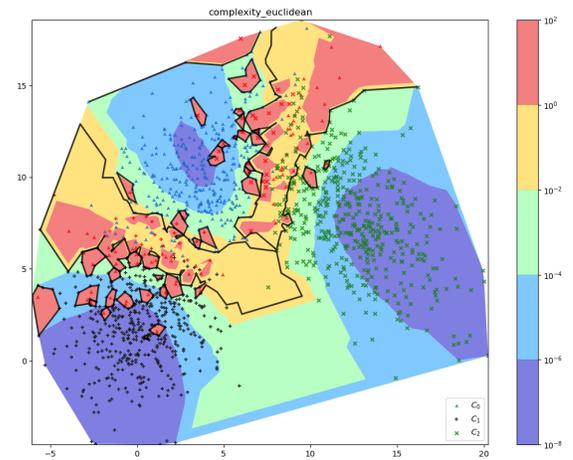


Figure 6: Three Gaussians with two distinct overlapping areas; background is Euclidean complexity.

Feature(s)	Slice	acc_{SVC}	Size	Rank
compl_mah	0.14 ~ 62.43	0.25	178	0.90
compl_euc	0.07 ~ 71.13	0.34	209	0.79
confidence	0.47 ~ 0.83	0.67	310	0.52
x_1, x_2	-6.04 ~ 8.62	0.76	279	0.51
	3.06 ~ 7.81			
compl_cos	1.08 ~ 1.60	0.66	255	0.50
x_1, x_2	5.90 ~ 10.95	0.76	243	0.49
	7.83 ~ 18.56			
x_1	5.90 ~ 9.45	0.74	207	0.48
x_2	12.14 ~ 18.56	0.78	165	0.45

Table 4: Slices for three Gaussians with multiple overlapping areas. Three Gaussians with two distinct overlapping areas; background is the model confidence. Top results of the data shown in Figures 5, 6 and 4

3.2 Real textual data

The experimental methodology performed in the real-life textual datasets is similar to the methodology used in the synthetic data except that here, the complexity is calculated on the test set while the class geometries are estimated on the train set. Also, as expected, the classifier was trained on the train set and evaluated on the test set. The split between train and test samples was performed randomly, maintaining 80%-20% train-test ratio.

We also experimented with following four textual datasets:

- ATIS (Hemphill, Godfrey, and Doddington 1990) is a standard benchmark data-set widely used as an intent classification. The test part has 800 records in 8 classes. The model we trained on it had 0.97 accuracy, meaning only 22 errors.
- WEBRR (Web Answer Passages; Keikha, Park, and Croft 2014) is a TREC GOV2 collection with 80 questions that serve as the labels and their answers, which serve as the samples; the answers were ranked and validated by users.
- NEWSGROUPS (Misra and Grover 2021; Misra 2018) is a collection of newsgroup documents. We used the 10 classes with the most documents. The model performed rather poorly with an accuracy of 0.57 (1060 errors).
- HWU64 (Liu et al. 2019) is Part of the DialoGLUE Benchmark containing popular personal assistant queries. The test set has 4108 records. The trained model achieve 0.80 accuracy (814 errors).

Table 5 provides a basic statistical overview on these dataset. Note that the normalized entropy of the samples’ distribution over the classes is used to estimate the degree of imbalance in the dataset. We formally define it as

$$\mathcal{E}(\mathcal{D}) = \frac{-\sum_{c \in \mathcal{Y}} \frac{|C(c)|}{|\mathcal{D}|} \cdot \log \frac{|C(c)|}{|\mathcal{D}|}}{\log(|\mathcal{Y}|)} \quad (9)$$

With the ATIS data-set the model confidence was by far superior in its estimate of how good it would do on the records. All other methods had similar results, with the

Dataset (\mathcal{D})	$ \mathcal{Y} $	$ \mathcal{D} $	$\mathcal{E}(\mathcal{D})$	MCS	acc_B
ATIS	8	4834	0.46	152	0.97
HWU64	60	20534	0.91	229	0.97
WEBRR	80	3298	0.88	25	0.57
NEWSGROUPS	41	50555	0.99	1250	0.8

Table 5: Dataset distribution statistics. In accordance with the notation in Section 2, $|\mathcal{Y}|$ denotes the number of classes, $|\mathcal{D}|$ denotes the total number of samples. $\mathcal{E}(\mathcal{D})$ is the the normalized entropy as defined in Equation 9, MCS is the median class size and acc_B denotes baseline accuracy according to the baseline classifier defined in Section 2.1

Mahalanobis method being considered somewhat better by FreaAI when considering the trade-off between size of slice and its performance.

Table 6 shows examples of high complexity and low complexity utterances from these data-sets. It can be easily seen that complex utterances tend to be longer and include words which by themselves could belong to either class. The complex utterances also tend to be more ambiguous in their wording.

To assess the affect of employing the correlation matrix instead of the covariance matrix when calculating the Mahalanobis distance, as suggested in Section 2, in Table 8 we report the ranking and the order that FreaAI yields for the two options for each of the datasets. We can see that on 3 of the 4 datasets the ranking of the Mahalanbis complexity measure is higher than when using the covariance matrix.

4 Related work

The geometrical properties underlying our complexity measures are based on GDA (Gaussian discriminant analysis(Ghojogh and Crowley 2019)) and thus underlie a wide variety of discrimination algorithms. This enables to use our measures as a black box, independent of the actual discrimination algorithm used to solve the classification problem. Our complexity methods are similar to those developed for the purpose of separability (Thornton 2008). The motivation for the separability measures is often feature selection (Mthembu and Marwala 2008; Guyon and Elisseeff 2003; Navot et al. 2005), dealing with the features to select to best separate the data, therefore achieving better discrimination among classes or clustering, and the evaluation of different clustering techniques (Peterson 2011; Guan and Loew 2021). Our motivation is different and is focused on automatically setting a baseline for testing ML discriminators. The different motivation also entails differences in implementation. For example, our work is efficiently different, as we consider the relationship between the different classes as our basis of complexity measure. Also, our method provides a score per record for the same computational cost, which is much more efficient than the common practice methods such as leave-one-out or Shapley values.

Dataset	$h(x, y)$	y	\hat{y}	Utterance
ATIS	3.68	ftime	flight	what are the departure times from detroit to westchester county
	3.15	flight	aircraft	show me the connecting flights between boston and denver and the types of aircraft used
	0.05	abbrv	abbrv	what does fare code qo mean
News	13.47	taste	good	These Heroic Food Trucks Are Coming To The Rescue Of California Fire Victims
	11.31	healthy	crime	Golden Gate Bridge Finally Getting A Suicide Barrier
	0.11	crime	crime	Two Police Officers Killed In Palm Springs, California Shooting
Web RR	23.62	rule	relation	If there are entrances in more than two directions, the union should be to the East.
	14.25	cult	church	In 1989,two members of a church, described by ATP as a doomsday religious cult...
	0.001	orange	orange	Tangerines

Table 6: Utterance examples: Examples of complex and simple utterances from the data-sets. Column $h(x, y)$ is the complexity measure based on the Mahalanobis distance function; y and \hat{y} are the utterance class and predicted class. Class name are shortened to fit the limited table space, especially for Web RR where these are phrases.

Dataset	acc	Metric	Slice	Slice Size	Rank	
			acc			
ATIS	0.97	<i>confidence</i>	0.03 ~ 0.47	0.15	26	0.86
		compl_mah	0.31 ~ 3.69	0.88	124	0.52
		compl_euc	2.23 ~ 4.54	0.82	93	0.49
		compl_cos	3.81 ~ 4.19	0.79	72	0.46
Web RR	0.97	confidence	0.02 ~ 0.17	0.69	52	0.89
		<i>compl_mah</i>	4.31 ~ 54.41	0.89	153	1.00
		compl_euc	7.93 ~ 20.16	0.86	121	0.93
		compl_cos	14.18 ~ 17.58	0.85	112	0.85
News	0.57	confidence	0.00 ~ 0.47	0.22	1294	0.57
		compl_mah	4.70 ~ 13.47	0.15	290	0.76
		<i>compl_euc</i>	4.46 ~ 8.05	0.12	382	0.89
		compl_cos	4.93 ~ 6.08	0.12	357	0.85
HWU64	0.80	<i>confidence</i>	0.00 ~ 0.35	0.31	1133	0.79
		compl_mah	1.35 ~ 119.6	0.59	1680	0.58
		compl_cos	13.63 ~ 17.67	0.57	1574	0.57
		compl_euc	6.05 ~ 29.01	0.55	1514	0.56

Table 7: Summary of the results for all textual data-sets. The table summarizes the results over all text data-sets. **Acc** is the accuracy of the model trained on the data-set, **Metric** is the information FreaAI used to search for slices, **Slice** is the range on the feature (except for confidence, the maximum is also the maximum in the data), **Slice Acc** is the accuracy of the records in the slice using the trained model, **Size** is the support or number of records in the slice and **Rank** is the FreaAI score for the slice (higher is more important). The metric which, according to FreaAI algorithm, provides the best information is marked in bold plus italic. Except for ATIS the complexity measures, especially the Mahalanobis one, are fairly close or better than the confidence of a *trained* model which means we can estimate before training a model areas which are hard to learn.

Dataset (\mathcal{D})	δ_M Rank	δ_M Order	$\delta_{\hat{M}}$ Order	$\delta_{\hat{M}}$ Rank
ATIS	0.54	2	0.88	2
HWU64	0.56	3	0.6	2
WEBRR	0.94	2	0.89	1
NEWSGROUPS	0.53	4	0.75	3

Table 8: Comparing FreaAI ranking score and order of the Mahalanobis-based complexity using the estimated covariance, δ_M , versus using the Pearson correlation matrix, $\delta_{\hat{M}}$.

5 Conclusions and discussion

We developed complexity measures and demonstrated their usage and usefulness for automatically setting a classification baseline as they highlight observations that are likely to be misclassified regardless of the discrimination algorithm used to classify them. Our method provides insights that are explainable via the complexity measure geometry, and does so at a linear calculation cost. We demoed our measures and methodology both over synthetic numeric data where we could easily validate the results, and over real natural language chatbot data. We believe that our complexity measures can have additional usages, such as assessing data quality and aiding in system design.

Our method may have several shortcomings including: (1) The complexity measure is calculated based on a text embedding and thus is limited by the embedding representation power. Therefore, one should select the embedding that best captures the important semantic aspects of the task. (2) The estimation of the class geometrical properties may not be accurate in small data-sets and thus one should take into account the size of the data-set when selecting the geometrical features to consider. (3) The complexity measure may return biased results in highly unbalanced data-sets. However, this limitation may be overcome using under-sampling techniques, similar to the methods employed in the field of active-learning to obtain a representative but smaller set of samples in each class. Alternatively, if possible, data augmentation techniques can be employed on the smaller

classes to obtain a more balanced data-set.

References

- Ackerman, S.; Raz, O.; and Zalmanovici, M. 2020. FreAI: Automated Extraction of Data Slices to Test Machine Learning Models. In Shehory, O.; Farchi, E.; and Barash, G., eds., *Engineering Dependable and Secure Machine Learning Systems*. Springer International Publishing.
- Anaby-Tavor, A.; Carmeli, B.; Goldbraich, E.; Kantor, A.; Kour, G.; Shlomov, S.; Tepper, N.; and Zwerdling, N. 2020. Do Not Have Enough Data? Deep Learning to the Rescue! *Proceedings of the AAAI Conference on Artificial Intelligence*, 34(05): 7383–7390.
- Barash, G.; Farchi, E.; Jayaraman, I.; Raz, O.; Tzoref-Brill, R.; and Zalmanovici, M. 2019. Bridging the Gap between ML Solutions and Their Business Requirements Using Feature Interactions. In *Proceedings of the 2019 27th ACM Joint Meeting on European Software Engineering Conference and Symposium on the Foundations of Software Engineering, ESEC/FSE 2019*, 1048–1058. New York, NY, USA: Association for Computing Machinery. ISBN 9781450355728.
- Galton, F. 1877. *Typical laws of heredity*. William Clowes and Sons.
- Ghojogh, B.; and Crowley, M. 2019. Linear and Quadratic Discriminant Analysis: Tutorial. arXiv:1906.02590.
- Guan, S.; and Loew, M. 2021. A distance-based separability measure for internal cluster validation. *arXiv*.
- Guyon, I.; and Elisseeff, A. 2003. An Introduction to Variable and Feature Selection. *J. Mach. Learn. Res.*, 3(null): 1157–1182.
- H., D.; C., B.; L., K.; A., S.; and V., V. 1996. Support vector regression machines. *Advances in Neural Information Processing Systems* 9.
- Hemphill, C. T.; Godfrey, J. J.; and Doddington, G. R. 1990. The ATIS Spoken Language Systems Pilot Corpus: Proceedings of a Workshop Held at Hidden Valley, Pennsylvania, June 24-27, 1990. www.kaggle.com/siddhadev/atis-dataset-from-ms-cntk.
- Keikha, M.; Park, J. H.; and Croft, W. B. 2014. Evaluating Answer Passages using Summarization Measures. *SIGIR*. <https://ciir.cs.umass.edu/downloads/WebAP/>.
- Ledoit, O.; and Wolf, M. 2004. A well-conditioned estimator for large-dimensional covariance matrices. *Journal of multivariate analysis*, 88(2): 365–411.
- Lee, K.; Lee, K.; Lee, H.; and Shin, J. 2018. A simple unified framework for detecting out-of-distribution samples and adversarial attacks. *Advances in neural information processing systems*, 31.
- Liu, X.; Eshghi, A.; Swietojanski, P.; and Rieser, V. 2019. Benchmarking natural language understanding services for building conversational agents. *arXiv preprint arXiv:1903.05566*.
- McLachlan, G. J. 1999. Mahalanobis distance. *Resonance*, 4(6): 20–26.
- Misra, R. 2018. News Category Dataset.
- Misra, R.; and Grover, J. 2021. *Sculpting Data for ML: The first act of Machine Learning*. ISBN 9798585463570.
- Mthembu, L.; and Marwala, T. 2008. A note on the separability index. arXiv:0812.1107.
- Navot, A.; Gilad-Bachrach, R.; Navot, Y.; and Tishby, N. 2005. Is Feature Selection Still Necessary? In *International Statistical and Optimization Perspectives Workshop "Subspace, Latent Structure and Feature Selection"*, 127–138. ISBN 978-3-540-34137-6.
- Peterson, A. D. 2011. *A separability index for clustering and classification problems with applications to cluster merging and systematic evaluation of clustering algorithms*. Ph.D. thesis, Iowa State University.
- Reimers, N.; and Gurevych, I. 2019. Sentence-BERT: Sentence Embeddings using Siamese BERT-Networks. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing*. Association for Computational Linguistics.
- Thornton, C. 2008. Separability is a Learner’s Best Friend. In *Proceedings of the Fourth Neural Computation and Psychology Workshop*. ISBN 978-3-540-76208-9.