



DEPARTMENT OF ENGINEERING MATHEMATICS

Improving Text Classifier Performance through Human-in-the-Loop: Enhancing Learning from Explanations

Xinyang Song

A dissertation submitted to the University of Bristol in accordance with the requirements of the degree
of Master of Science in the Faculty of Engineering.

Sunday 20th August, 2023

Supervisor: Dr. Edwin Simpson

Declaration

This dissertation is submitted to the University of Bristol in accordance with the requirements of the degree of MSc in the Faculty of Engineering. It has not been submitted for any other degree or diploma of any examining body. Except where specifically acknowledged, it is all the work of the Author.

Xinyang Song, Sunday 20th August, 2023

Contents

1	Introduction	1
1.1	Background	1
1.2	Motivation	1
1.3	Experimental Methodology	2
2	Background	5
2.1	ExpBERT	5
2.2	HITL	6
2.3	Active Learning	6
2.4	Small-Text Library	10
2.5	GPT Annotator	11
2.6	Bayesian Neural Networks in Active Learning	11
3	Design	13
3.1	Structure Overview	13
3.2	Pre-training	14
3.3	Classifier Model	15
3.4	Active Learning Strategy	16
3.5	Annotation Process	18
4	Implementation	21
4.1	Experimental Environment	21
4.2	Dataset	21
4.3	Performance Metrics	23
4.4	Stopping Strategy	23
4.5	Execution Module	24
5	Critical Evaluation	27
5.1	Expectations and Evaluation Framework	27
5.2	Expectation One	28
5.3	Expectation Two	28
5.4	Expectation Three	29
5.5	Expectation Four	30
5.6	Expectation Five	31
5.7	Conclusion	33
6	Conclusion	35
6.1	Project Status	35
6.2	Limitations	36
6.3	Future Work	36
A	Github Repository	45

List of Figures

2.1	ExpBERT to produce representations that form the input to the classifier	5
2.2	Pool-based active learning process	7
2.3	Traditional Neural Networks and Bayesian Neural Networks	11
3.1	The overall experimental framework	13
3.2	Overview of pre-training process	14
3.3	Standard Neural Network and Dropout Neural Network.	16
4.1	Average F1 score of full-data and random sampling strategies under validation dataset . .	24
5.1	Random sampling only annotates labels	28
5.2	Average accuracy score for test set	30
5.3	Average F1 score for test set	30
5.4	Average accuracy score for validation set	30
5.5	Average F1 score for validation set	30
5.6	Average accuracy score in test set	32
5.7	Average F1 score in test set	32

List of Tables

3.1	Results of the OpenAI model evaluation.	19
4.1	Data distribution in CrisisNLP.	22
4.2	Distribution of datasets	22
5.1	Evaluation structure.	27
5.2	Results of evaluation ID 2 and evaluation ID 3	29
5.3	Performance of different strategies through iterations	29
5.4	Performance of different strategies in the last iteration	31
5.5	Results of evaluation ID 5 and evaluation ID 6	32

List of Algorithms

3.1 BALD using MCD	17
------------------------------	----

Abstract

Text classification utilises **Natural Language Processing (NLP)** techniques to analyse pre-trained texts and assign them appropriate labels. Especially in crisis situations like floods, earthquakes, etc., text classifiers are crucial in identifying key information and effectively forwarding it to relevant agencies. However, the efficacy of text classification largely depends on abundant rich training data, which may be difficult to obtain in many scenarios [45]. Relying on a vast amount of non-representative annotated data can lead to delays in project commencement and may impact the model’s accuracy. Specifically, in emergencies, identifying action-related information (such as casualties or missing persons) becomes particularly challenging.

To address these issues, this paper introduces a Human-in-the-Loop (HITL) system integrated into the classifier’s training process. Coupled with the representational engineering technique of natural language explanations—ExpBERT, a BERT fine-tuned on the MultiNLI natural language inference dataset is utilised to learn from explanations [52]. This optimised embedding representation is used as input to the neural network classifier, further enhancing performance. The core task of the paper revolves around improving the performance of the text classifier based on ExpBERT by combining a limited but representative set of samples with labels and explanations provided by either humans or OpenAI [5].

This paper delves deeply into the application of active learning in text classification tasks based on ExpBERT. During the iterative process, acquisition functions based on uncertainty and diversity are used to select representative unlabeled instances for annotators to process. The original small sample and newly annotated data are used for model retraining. The results indicate that this method can achieve performance similar to models trained on vast datasets within a few iterations. To substantiate this, the model’s performance using different active learning strategies is compared on the CrisisNLP dataset, and the effects of an active learning system combined with Monte Carlo Dropout (MCD) are also evaluated.

The main conclusions of this paper are:

- The Human-in-the-loop framework combined with a text classifier enables the model to achieve, or even surpass, the performance of models trained with the full dataset, using only a small amount of annotated data and undergoing a limited number of iterations.
- Compared to merely annotating labels, low-quality explanations, or a limited number of explanations, annotating a certain amount of high-quality explanations can significantly boost the model’s performance.
- In terms of acquisition functions, semantic diversity-based sampling and Bayesian Active Learning by Disagreement (BALD) strategies combined with Monte Carlo Dropout (MCD) can allow the model to achieve higher average performance post-training.



Supporting Technologies

- I used Python as the development language.
- I used the *Transformers* library to introduce pre-trained models and *Pytorch* for implementing neural networks.
- I used OpenAI models as one of the annotation methods.
- I used L^AT_EX to format my thesis, via the online service *Overleaf*.

Notation and Acronyms

NLP	:	Natural Language Processing
BERT	:	Bidirectional Encoder Representations from Transformers
ExpBERT	:	Representation Engineering with Natural Language Explanations
NLI	:	Natural Language Inference
HITL	:	Human-in-the-loop
AL	:	Active Learning
BALD	:	Bayesian Active Learning by Disagreement
MCD	:	Monte Carlo Dropout
GPT	:	Generative Pre-trained Transformer
NN	:	Neural networks

Acknowledgements

Xinyang Song would like to extend her profound gratitude to her supervisor, Dr. Edwin Simpson, for sharing literature and implementation examples related to ExpBERT and active learning, offering countless invaluable assistance. His professional guidance and enthusiastic support helped me overcome numerous challenges, complete my dissertation, and achieve meaningful insights.

Additionally, I'd like to thank the University of Bristol and my peers for their resources and emotional support during this time. My gratitude goes out to all of you.

Chapter 1

Introduction

This **chapter** begins by introducing the research background of the text classifier based on ExpBERT integrated with Human-in-the-loop (Section 1.1). Based on this background, the motivation for the experiment was formed in Section 1.2, emphasising the significance of this research. Furthermore, a general overview of the work will be provided, contrasting with the limitations of traditional methods. Finally, the primary objectives and challenges are briefly summarised (Section 1.3).

1.1 Background

Emergency response systems for social platforms mainly focus on developing better text classification algorithms to learn from limited data. However, obtaining valuable annotated datasets can prove to be challenging [60]. The text classifier based on ExpBERT receives classification labels and corresponding explanations. These explanations detail which keywords led to a particular classification. ExpBERT incorporates this information into the model’s training, allowing it to learn from deeper semantic insights to enhance its generalisation capabilities. Therefore, annotating with representative data and understanding the key phrases is vital for guiding the model towards accurate classification.

However, the data throughput of typical social platforms does not permit the annotation of tens of thousands of entries. Consequently, various forms of active learning that can sample high-information-value data at low costs have found extensive application in classification projects [61, 76]. Active learning emerges as an effective solution to these challenges, selecting a limited number of highly informative unannotated samples for human experts in the loop to annotate [69, 19]. The newly data, once effectively annotated, dictates the performance of the model in the next iteration. Thus, employing different acquisition functions to query the most information-rich new instances is perhaps the most popular approach in active learning. Acquisition functions have naturally become the focal point of research in the domain of active learning.

Additionally, neural network-based text classifiers are often ill-suited for early uncertainty sampling [67, 64]. This is because the weight parameters of the neural network are fixed values, causing the model to be overly confident in its predictions, both correct and incorrect. However, by employing dropout, it’s possible to introduce a degree of uncertainty into the model. During training, the dropout layer randomly “turns off” a subset of input units. This randomness ensures that network weights are no longer fixed values, thus simulating weight uncertainty.

1.2 Motivation

In light of the aforementioned background, this section emphasises the research motivation from three aspects: accuracy, representativeness, and robustness.

1.2.1 Accuracy

A unified feature of social networks, exemplified by Twitter, is the vast amount of data combined with internet-specific language syntax. The ultimate goal is to precisely classify urgent messages when capturing emergent needs and to dispatch relevant departments to address those needs. However, traditional text classification systems that achieve performance improvements are primarily done through extensive supervised learning, which is costly and slow to respond. Texts with varying amounts of information are randomly distributed, and the difference in information content can significantly impact model accuracy. Consequently, the results of capturing and identifying crucial emergency information are often subpar. Therefore, employing active learning models to utilise a small amount of data and enhance accuracy is urgently needed.

1.2.2 Representativeness

The ability to query rich feature representations from a large pool of unannotated data can mitigate much of the bias in multi-class active learning problems for emergency scenarios. By actively selecting samples for annotation, the model can choose the most valuable and representative samples in each iteration to enhance its performance. Therefore, selecting an effective acquisition function to seek out representative data has the most profound impact on performance and offers the most significant room for improvement.

1.2.3 Robustness

Original models often exhibit an overconfidence bias when handling text classification tasks [36]. When encountering unfamiliar content, they make incorrect and unreliable predictions. In many practical applications, such as disaster response and medical diagnostics, erroneous predictions can have severe consequences. Thus, introducing a certain level of uncertainty to the model to curb its overconfidence has become an essential research topic. On the other hand, we can enhance the robustness of the model without introducing excessive computational overhead. This holds significant value for neural networks in scenarios that demand highly accurate and reliable predictions.

1.3 Experimental Methodology

Considering the background and motivations, this paper eventually designs a pool-based Human-in-the-loop active learning system for text classification. Unlike traditional text classification models, we train with a small amount of **annotated data** and then use the trained model in combination with various acquisition functions to select instances. Subsequently, preset annotations or human or OpenAI annotators observe and analyse the key features of the extracted instances in the loop, providing corresponding labels and adding a certain amount of explanations. These newly annotated instances are incorporated into the original data, and the above cycle is repeated until the performance reaches that of the full data (with the training set ratio close to 1, annotated with labels and default explanations).

The implementation will utilise uncertainty sampling based on Least Confidence, sentiment diversity sampling, and random sampling as a baseline to examine the performance under various active learning strategies. Detailed explanations of the acquisition functions will be provided in Chapter 2. Furthermore, in every iteration, the concatenated explanatory texts will be processed using a pre-trained model based on ExpBERT to generate deep semantic representations. This embedded vector will then be input into the neural network text classifier.

Lastly, structural modifications will be made to traditional neural **network(NN)** models. Using the dropout mechanism, the model becomes more robust to minor variations in input, thereby enhancing its generalisation performance. Additionally, dropout can address the model's overconfidence issue by introducing noise and randomness, ensuring that model predictions aren't overly deterministic [59]. While dropout doesn't directly quantify uncertainty, the introduced randomness and noise can bolster the model's resilience and, to some extent, alleviate its overconfidence. Given the high computational cost of entropy calculations in the original Bayesian Active Learning by Disagreement (BALD) algorithm, we've

refined the BALD algorithm. Utilising Monte Carlo Dropout (MCD), we employ Least Confidence to select samples where the model’s predictive probability is most uniform (i.e., without a particularly high predictive probability for any category).

1.3.1 Objectives

Based on the above experimental descriptions, the objectives of this paper can be summarised as follows:

- Develop a pool-based Human-in-the-loop framework applied to the ExpBERT-based text classifier and investigate its effectiveness.
- Construct an annotator simulation process and design stopping criteria (e.g., terminate when reaching the performance of the full data model).
- Establish a baseline (performance of the full-data model without active learning or random sampling) and study whether active learning can enhance the performance of the text classification model. Compare the effects of different acquisition functions.
- Explore the impact of the quality and quantity of explanations on this system. For instance, providing explanations with noise or increasing the number of explanations added in each iteration.

1.3.2 Challenges

The main challenges of this project lie in the significant differences that may exist between various active learning strategies and how to improve the original initialisation framework of the multi-classifier based on ExpBERT and the annotation process to reduce computational costs.

- It’s worth noting that there are multiple acquisition functions in active learning text classification systems. The application of different algorithms in active learning strategies significantly affects the performance of the active learning framework. The choice of strategy needs to consider both time complexity and representativeness.
- Secondly, neural networks with dropout mechanisms might increase computational costs during training algorithm execution and when using MCD for active learning. Thus, how to reduce computational costs in each iteration to improve response efficiency poses a challenge.
- Finally, regarding the framework design, determining the appropriate sampling ratio, the number of iterations, and the number of explanations provided in each iteration requires further study to achieve satisfactory performance.

Chapter 2

Background

Chapter 2 introduces the relevant techniques centred around enhancing the performance of the text classifier based on ExpBERT through active learning. The ExpBERT mechanism used in the experiment is introduced (Section 2.1). ExpBERT will be applied during dataset pre-training. The active learning-based “Human-in-the-loop” (Section 2.2) and its main acquisition functions (Sections 2.3 and 2.4) are the focal implementation parts of this experiment, which can rapidly boost performance through this technique. Section 2.5 describes the multi-model integration method, while Section 2.6 discusses using OpenAI models to simulate the annotator role within the Human-in-the-loop. Lastly, the feasibility of applying Bayesian theory to active learning is concerned.

2.1 ExpBERT

Representation Engineering with Natural Language Explanations (ExpBERT) proposes a method to enhance the control over the inductive biases of a model and maintains the expressive power of the representation [52]. It combines fixed explanations provided by humans or, in this experiment, annotators with tweets to learn from these explanations and improve the model’s performance. Figure 2.1 intuitively shows how to combine samples with explanations with BERT fine-tuned on MultiNLI [77]. Explanations play a key role here. The quality of the explanations has a significant impact on the performance of ExpBERT. By using high-quality explanations, the model’s learning can be guided.

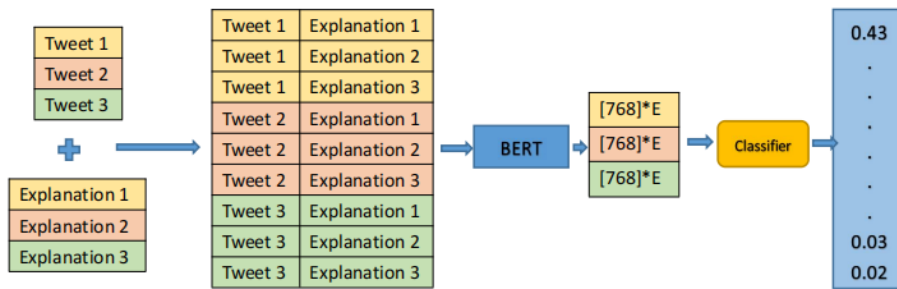


Figure 2.1: ExpBERT to produce representations that form the input to the classifier

First, each tweet is fully connected with a set of pre-written explanations to generate a 3×3 set. These explanations are unrelated to specific tweets; each tweet will be connected with the same number of explanations. The preprocessed text and explanations are input into BERT to generate features that “interpret” each explanation [15]. The classifier can be trained to classify the representations of the tweets and explanations.

After being fine-tuned on the MultiNLI dataset, the BERT model will generate a feature vector for each input sample, representing the entire input of a length of 768. The feature vectors of tweets and explanations are then connected to form an embedding with a size of $768 * E$, where E is the number

of explanations. This will be used as the input data for training and predicting in the classifier model. This paper will use this model as the basic pre-trained model. At the same time, combined with Natural Language Inference (NLI) technique can be used to reduce the number of embeddings (the size will be $768 + (3 * E)$) and integrate and initialise them into the HITL system [42]. Because this model can handle instance vectors with explanations, the optimised embedding representation as input to the classification model can improve the performance of the model.

2.2 HITL

Traditional natural language processing workflows are not designed to fully leverage human feedback. In contrast, Human-in-the-Loop (HITL) serves as an essential component of interactive systems, revealing potential model flaws through simulating human roles within the loop. These flaws might not be readily apparent before real-world testing [74]. Godbole et al. (2004) [23] expanded a text classifier employing Support Vector Machine (SVM) active learning. Ingeniously, they integrated human input on aspects such as feature engineering, terminology selection, and document labelling into the model, allowing statistically sound judgments to be made. This innovative mode of interaction between humans and machine learning algorithms is termed “Human-in-the-Loop (HITL) machine learning” [49]. In this approach, various HITL machine-learning schemes can be determined based on different needs and contexts, driven by different types of collaborations between humans and machines [50].

Active Learning (AL) [69]: The essence of active learning lies in the system’s control over the model’s learning process. Although humans serve as intermediaries to annotate unlabeled data, they cannot choose unlabeled data based on preference. AL will be applied as an optimisation framework in this paper, with a detailed exploration of its application in Section 2.3.

Interactive Machine Learning (IML) [2] is a method where humans maintain close interaction with the system. Within this context, there are notable distinctions between Active Learning (AL) and Interactive Machine Learning (IML). Dudley and Kristensson (2018) [17] highlighted that although both AL and IML prioritise the user’s selection of new data points for labelling, they are driven by different motivations. Specifically, in AL, the selection is propelled by the model itself, while in IML, the choice is driven by the user. Given that IML has evolved from the foundations of AL, they share several common drawbacks. However, IML also faces unique challenges, especially those related to Human-Computer Interaction (HCI) complexities. These distinct challenges necessitate more specialised research to address, as indicated by Michael et al. (2020) [47].

Machine Teaching (MT) [63] is an approach where machine learning models are trained under the guidance of human educators. More precisely, human teachers define the knowledge they aim to convey to the model, underscoring their active involvement and direction throughout the learning process. Devdize et al. (2020) [14] pointed out that, compared to active learning, MT heavily leans on the teacher’s expertise and offers less flexibility in sample selection and handling intricate tasks. Hence, selecting the right learning method becomes paramount based on the specific requirements.

2.3 Active Learning

Within the context of HITL, active learning systems, as one of the most popular learning approaches, aim to overcome the labelling bottleneck by posing questions to unlabeled instances and having them labelled by experts, such as human annotators [1, 69]. In essence, active learning seeks to identify the most informative unlabeled raw data and present it to annotators for labelling. This process closely mirrors real-world scenarios of extracting information from source data. Annotators label this source data, and the labelled instances are incorporated into the model’s training process. In this way, active learning manages to achieve performance comparable to using the full dataset but with fewer labelled data, thereby addressing the labelling bottleneck. The active learning approach proves especially pertinent given this project’s vast number of instances, coupled with significant noise, cumbersome full-scale training tasks, and suboptimal accuracy. Thus, active learning is apt for this project, focusing on instances that might enhance accuracy, minimise irrelevant instances, and boost the learning efficiency and precision of the

emergency response system.

2.3.1 The AL Process and Scenarios

The active learning process is illustrated in Figure 2.2, which depicts the process of pool-based active learning. The machine learning model is initially trained using the labelled instance set L . Afterwards, the model gives the score to the unlabeled dataset and selects representative unlabeled samples based on the acquisition function, presenting them to human annotators. These annotators label the chosen samples, which are then removed from the unlabeled sample set U and added to the labelled set L . Through iterative execution of the steps above, the quantity of labelled samples grows incrementally, enhancing the model’s performance. The active learning process will conclude once the stopping criterion is met, such as achieving performance metrics on the entirety of the data set.

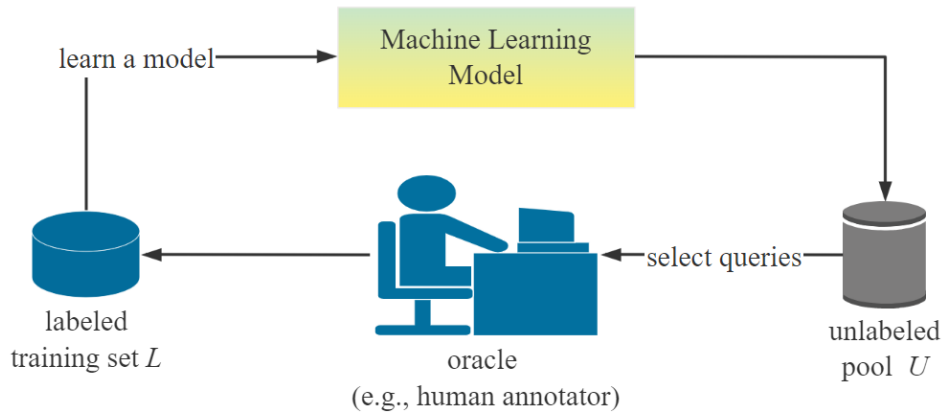


Figure 2.2: Pool-based active learning process

Based on the classification by Settles (2009) [69], active learning primarily encompasses three scenario settings: Membership Query Synthesis [3, 34], Stream-based Selective Sampling [11], and Pool-based active learning [40].

Membership Query Synthesis: In this active learning scenario, the strategy is generated by the model itself. The machine learning model can request labels for any unlabeled instance that hasn’t been sampled from the underlying natural distribution [69]. This query approach is typically effective for more deterministic problem domains, such as tasks involving regression predictions on absolute coordinates. It can interpret straightforward data distributions and construct viable data for human annotation. However, as evidenced in the study by Baum and Lang (1992) [37], in complex tasks like natural language processing, the model might produce text strings that are hard to comprehend, rendering human judgment on these jumbled texts challenging. Within the context of deep active learning, one can address the scenarios of Membership Query Synthesis through Generative Adversarial Networks (GANs) for data augmentation, as GANs are capable of generating instances with a high degree of plausibility [25].

Stream-based Selective Sampling: The distinction of this setup in comparison to Membership Query Synthesis lies in its ability to operate irrespective of the input distribution. Stream-based Selective Sampling entails drawing a single unlabeled instance from the actual distribution at a time. The model then determines whether to request the label of the instance through “information measures” or “query strategies”, essentially acting as a kind of biased random sampling [11].

Pool-based active learning: As a commonly employed setting in active learning, the difference between pool-based active learning and stream-based Selective Sampling is that the former employs a greedy mechanism before selecting the best query, comparing the entire dataset. In contrast, stream-based sampling assesses data as it comes in on an individual basis. However, dealing with vast datasets could result in lengthy computational times due to the need to evaluate all the data. Nonetheless, this method is

particularly suitable for scenarios similar to this project where manual annotation costs are significant. By choosing the most valuable data for annotation, it aims to maximise the return on the annotation investment.

2.3.2 Acquisition Functions

This section considers various acquisition functions in active learning. In the design of the paper (Chapter 3), we will utilise a text multi-classification model based on ExpBERT combined with different acquisition functions to enhance performance.

Random Sampling

Random sampling involves selecting instances independent of predictions, data, or model considerations. It's commonly used as a task benchmark since it might overlook potential informative instances, reducing learning efficiency. In this context, random sampling will serve as a baseline to contrast with the more sophisticated strategies discussed below, especially when the annotation pool becomes overly extensive [68].

Prediction Uncertainty Sampling

Uncertainty sampling is a specific active learning strategy that queries instances most difficult for the model to classify. By annotating these challenging instances, the model's accuracy is enhanced. In probabilistic models for binary classification, such instances have posterior positive probabilities close to 0.5 [40]. However, for more complex multi-label classification tasks, entropy-based methods are used. The more uniform the probability distribution, the greater the entropy. As entropy increases, the uncertainty or informational content of the random variable becomes higher. Conversely, when probabilities are concentrated on a few data points, the uncertainty is lower and the informational content is smaller.

$$x_{\text{ENT}}^* = \arg \max_x - \sum_i P(y_i | x; \theta) \log P(y_i | x; \theta) \quad (2.1)$$

$P(y_i)$ denotes the probability distribution at classification i .

In the domain of text classification, an alternative method often used to measure uncertainty is the “Least Confidence(LC)” [30]:

$$x_{LC}^* = \arg \min_x P(y^* | x; \theta), y^* = \arg \max_y P(y | x; \theta) \quad (2.2)$$

y^* represents the label of the class with the highest likelihood. This method is equivalent to the effectiveness of entropy-based algorithms for binary text classification. It focuses on samples for which the model predicts the highest probability but with low confidence. Specifically, the method involves annotating the samples with the smallest maximum probability. This more comprehensive consideration will be adopted in this paper.

In addition to the two commonly used measurement methods mentioned above, Munro [49] introduced the margin of confidence and ratio of confidence. The former refers to the difference between the two most confident predictions, while the latter is the ratio between these two predictions.

$$\text{margin_conf} = 1 - (P(y_1^* | x) - P(y_2^* | x)), \quad \text{ration_conf} = \frac{P(y_2^* | x)}{P(y_1^* | x)} \quad (2.3)$$

y_1^* is the most confident, y_2^* is the second most confident.

Query-By-Committee

Compared to uncertainty sampling, which evaluates the judgement of a single model, the Query-By-Committee (QBC) approach employs multiple models to act as “members” of a committee [70]. These members are trained on the available labelled data using the Gibbs algorithm. The committee then randomly selects a hypothesis consistent with the current labelled data. The subsequent query is picked based on the instance where the committee members exhibit the most significant disagreement [9]. Hence,

while uncertainty sampling assesses the discernment of an individual member (model), QBC often incorporates strategies such as using vote entropy to select instances that are challenging for models to agree upon [11]. And choosing samples with a high average Kullback-Leibler (KL) divergence [44].

Vote entropy: For text classification tasks, each committee member (text classifiers) can vote on unlabeled samples.

$$H(V) = - \sum_{i=1}^n \frac{Vote(y_i)}{C} \log \frac{Vote(y_i)}{C} \quad (2.4)$$

$Vote(y_i)$ represents the votes received for class y_i , and the total number of votes is C , which is also the total number of classifiers. The larger the vote entropy, the more controversial the sample is.

Average KL Divergence: The KL divergence increases as the difference between two distributions grows. The formula is expressed as:

$$x_{KL}^* = \underset{x}{argmax} \frac{1}{C} \sum_{c=1}^C D_{KL}(P_c || \bar{P}) \quad (2.5)$$

Where,

$$D_{KL}(P_c || \bar{P}) = \sum_i P_c(i) \log \frac{P_c(i)}{\bar{P}(i)} \quad (2.6)$$

P_c is the prediction probability distribution of the i th committee member, and \bar{P} is the average prediction probability distribution of all members.

It is noteworthy that Query-By-Committee (QBC) requires simultaneous operation of multiple models. Hence, one drawback of methods based on QBC is their slower speed, which is a factor of the committee's size [7]. Speed is paramount for active learning applications, but the trade-off with QBC is that annotators experience prolonged waiting times while the system identifies instances using QBC.

Semantic-based Diversity Sampling

Peng, Hao, et al. [58] introduced a semantic-based diversity sampling method that can be applied in text classification combined with active learning in this experiment. Distinct from the uncertainty sampling approach, the semantic-based diversity sampling method eliminates text sample redundancy semantically using the *Euclidean distance* [62]. This distance measure is frequently used across various domains, including data clustering and nearest neighbour search in embedding spaces. In machine learning, calculating a Euclidean distance from an embedding to a cluster centre can be achieved by computing a distance matrix for every embedding to every cluster centre, where “ $distances[i, j]$ ” represents the Euclidean distance from the i th embedding to the j th cluster centre [6].

To ensure that richer and less redundant samples are provided to the model (learner) in subsequent processes, this abstraction method employs the greedy k-centre clustering algorithm by Sener and Savaravarese (2017) [68]. The dataset D contains $n \times m$ unlabelled texts and divides D into n batches, with each sample set containing m instances. It is the result of encoding the dataset. First, select α vectors from Y_i^* to initialise the clusters c_i^0 . Here examples will be considered as cluster centres. Then the k-centred algorithm searches u_i^n from c_i^0 , which is a set that includes members that are not in Y_i^* . It is the furthest from the centre of all clusters. The algorithm chosen is formulated as follows:

$$\begin{aligned} u_i^0 &= \arg \max_{y_i^k \in \tilde{c}_i^0} \min_{y_i^j \in c_i^0} \|y_i^k - y_i^j\|_2^2 \\ u_i^1 &= \arg \max_{y_i^k \in \tilde{c}_i^1} \min_{y_i^j \in c_i^1} \|y_i^k - y_i^j\|_2^2 \end{aligned} \quad (2.7)$$

Where:

$$\begin{aligned} \tilde{c}_i^0 &= Y_i^* \setminus c_i^0 \\ c_i^1 &= c_i^0 \cup u_0 \end{aligned} \quad (2.8)$$

This is followed by updating the existing clusters to c_i^1 after a loop execution and merging the output into P . All text instances in P converge into a core set that best represents and generalises the dataset D in the semantic space.

Bayesian Active Learning by Disagreement

Bayesian Active Learning by Disagreement (BALD) is an active learning strategy suitable for text classification problems, combining Bayesian inference with uncertainty measures to select information-rich samples. In BALD, uncertainty is measured by calculating the inconsistency in predictions for each sample's class probability distribution, as proposed by Houlisby, N. (2011) [29]. The learner (model) aims to maximise the uncertainty of the model parameters through the input x . $H(y | x, D)$ represents the uncertainty of the target variable. The second term of the equation represents the expected uncertainty (entropy) of $H[y | x, \theta]$, given that the parameter θ follows the posterior probability distribution $p(\theta | D)$ based on the training dataset D , measuring the average uncertainty.

$$x^* = \arg \max_x H(y | x, D) - E_{\theta \sim p(\theta | D)} [H[y | x, \theta]] \quad (2.9)$$

To better calculate the inconsistency between predictions, Gal et al. proposed a specific sampling function in 2017 that utilises the Monte Carlo Dropout (MCD) method [21]. MCD is a regularisation technique during training, which, by randomly deactivating certain neurons, simulates the posterior distribution of Bayesian networks [53]. This allows for selecting samples with the highest uncertainty for labelling by making multiple predictions for each sample and calculating the inconsistency between the model's predictions.

At its core, MCD is a regularisation method designed to prevent overfitting by probabilistically deactivating a subset of neurons. During the inference process, different results are produced during the prediction phase by inputting into the neural network a finite number of times (T times) and using dropout [75]. Consequently, after quantifying the uncertainty from these varying prediction results, they can be used to estimate the model's uncertainty about unlabeled samples. This is then applied in the Bayesian Active Learning by Disagreement (BALD) sampling strategy to select the most uncertain unlabeled data. The inference results of MCD are derived in the following manner:

$$\bar{y}(x) \approx \frac{1}{N_{mc}} \sum_{N_{mc}} y_{drop}(x) \quad (2.10)$$

Wherein, y_{drop} represents the varying outputs from the dropout-enabled network, x is the input to the network, and N_{mc} is the number of samples required to obtain the distribution. Ultimately, BALD selects the samples with the highest informational gain by comparing each sample's BALD value. These samples are then labelled and added to the annotated training set. Moreover, BALD can be paired with the previously mentioned *Least Confidence*, replacing the original entropy computation, to select samples for which the model prediction probabilities are most uniform (i.e., samples without a particularly high prediction probability for any category).

2.4 Small-Text Library

The text classifier in this project tends to focus on a single model, potentially overlooking the applicability of other feasible models. However, the time-consuming of switching models and active learning strategies, coupled with code redundancy, would significantly impact the progress of the experiment. The Small-Text Library integrates commonly used libraries like scikit-learn, transformers, and PyTorch that can be applied within the Python environment [66].

The pool-based active learning framework for text classification bridges the interfaces of active learning strategies, classifiers, and stopping criteria. It provides an advanced active learning framework for text classification tasks and offers a range of classifier and acquisition function components. These can be combined in various permutations for experimental and applied active learning tasks, facilitating the implementation of active learning in the Python ecosystem. Compared to the popular ModAL library

[12], Small-Text offers more flexible customisation services, with the former focusing on model ensembles and strategy choices. Nonetheless, this all-encompassing approach might not necessarily enhance the performance of the ExpBERT-based project, as performance discrepancies can vary based on circumstances.

2.5 GPT Annotator

GPT (Generative Pre-trained Transformer) is a pre-trained language model developed by the OpenAI team in 2018 [4]. Its key algorithm is the Transformer, a deep neural network structure based on the self-attention mechanism, boasting potent sequence modelling and representation learning capabilities [41]. This model can analyse and generate natural language text through pre-training and fine-tuning, proving beneficial in multiple scenarios such as automatic answering, intelligent customer support, and language translation.

As a result, for a time-saving and efficient text classification system, this project will employ this model as an active learning annotator based on ExpBERT. The primary task is to generate appropriate explanations for ExpBERT and provide them to the model to enhance classification performance. In this configuration, active learning generates human-readable explanations for input text by leveraging GPT's robust generative model [72].

A common approach involves using GPT to find keywords from the sampled text of a specific category, combine the keywords to form explanations, and then choose one or more to integrate into the original explanations. The combined text is then passed to the ExpBERT model. By analysing the sampled text in this manner, one can overcome deficiencies in human memory and the inadequacy of semantic similarity computation, ultimately identifying crucial feature words and phrases that best assist the model in understanding the intrinsic structure of the data.

2.6 Bayesian Neural Networks in Active Learning

In modern machine learning and deep learning applications, Bayesian Neural Networks (BNNs) have emerged as critical components due to their flexible inference mechanisms. BNNs offer a rich and adaptable modelling framework by introducing uncertainties in weights and predictions [43]. This not only bolsters the model's reliability but also enhances its ability to generalise, adapt to noise, and handle outliers. Unlike traditional neural networks [38], where weights are fixed, for BNNs, they are variable, as shown in Figure 2.3. The idea behind the model is to combine the strengths of neural networks with probabilistic modelling, providing probabilistic assurance for predictions [51].

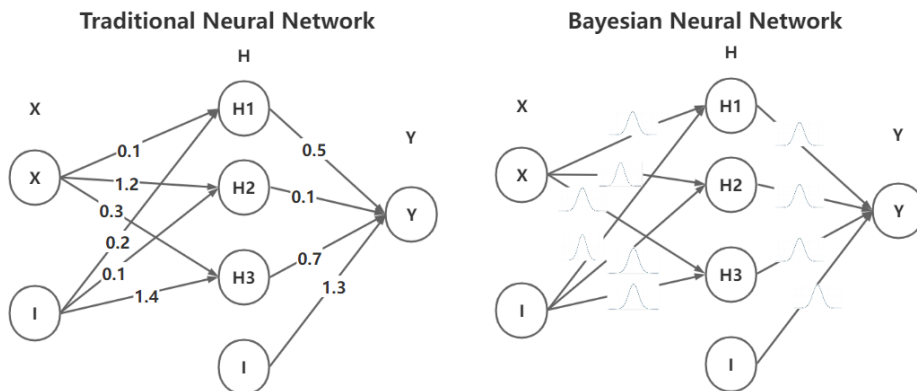


Figure 2.3: Traditional Neural Networks and Bayesian Neural Networks

The weights in this context are inferred based on the observations we have, presenting an inverse probability problem solved through Bayes' theorem [22]. As shown in the equation, the weight ω represents latent variables for which the true distribution is unseen. Bayesian inference allows us to obtain

the distribution of model parameters conditioned on the observed data, $p(\omega \mid D)$, known as the posterior distribution.

$$\pi(\omega \mid D) = \frac{p(\omega)p(D \mid \omega)}{\int p(\omega)p(D \mid \omega)d\omega} = \frac{p(\omega)p(D \mid \omega)}{p(D)} \quad (2.11)$$

Additionally, the likelihood function $p(D \mid \omega)$ in multi-class problems combines the neural network's predictions (after a Log Softmax transformation) with the actual observed class labels to quantify their consistency. The posterior distribution can be computed using Bayesian theory after determining the likelihood and the prior distribution $p(\omega)$.

After obtaining the posterior distribution, marginalisation methods will be used for prediction:

$$\mathbb{E}[f] = \int f(\omega)\pi(\omega \mid D)d\omega \quad (2.12)$$

Bayesian Neural Networks provide a way to quantify uncertainty, enabling active learning systems to identify instances on which the model is least confident, subsequently selecting these instances for annotation. However, integrating the model parameters often requires computationally expensive operations, such as Markov Chain Monte Carlo (MCMC) sampling [27]. This can become problematic when dealing with large-scale datasets, especially in the context of active learning, which typically requires multiple iterations and model updates. Therefore, the actual implementation of BNNs may prolong the initialisation time of emergency systems [54].

Chapter 3

Design

This chapter will provide an overview of the active learning framework in Section 3.1, followed by a detailed breakdown and design according to the sequence of tasks in active learning. Starting with pre-training the data in Section 3.2, we then discuss the design of classifier models that acquire embedded representations in Section 3.3. Section 3.4 presents the design of active learning strategies from multiple perspectives using the trained classifier model and analyses the implementation and effects of different acquisition functions. Finally, the design of the annotator’s process functionality is covered in Section 3.5.

3.1 Structure Overview

The overall experimental framework is built upon the pool-based active learning technique discussed in Section 2.3.1. The system samples from a fixed pool of unannotated data using a trained model for sample assessment.

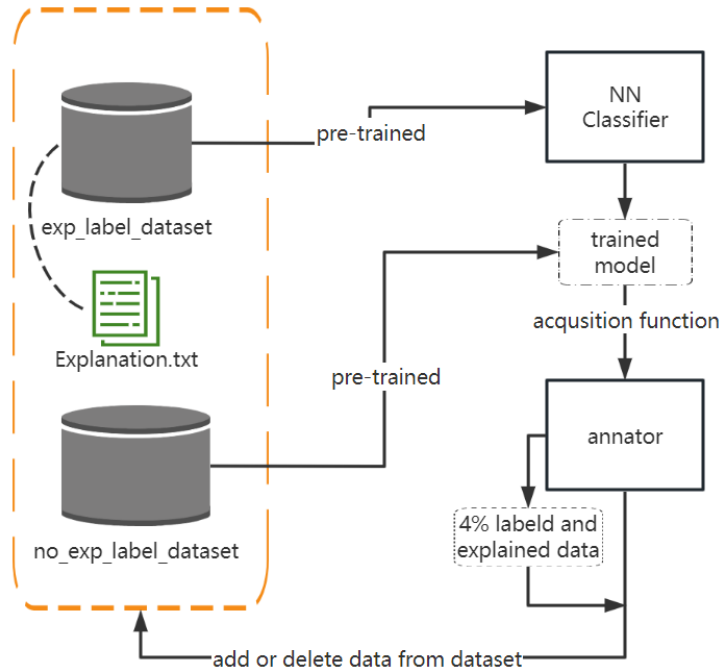


Figure 3.1: The overall experimental framework

As depicted in Figure 3.1, in contrast to traditional active learning methods, this study revises the dataset requirements, introducing two distinct datasets: **one annotated with labels and explana-**

tions (exp_label_dataset) and **another neither labelled nor explained** (no_exp_label_dataset). The ‘exp_label_dataset’ encompasses data for training, validation, and testing. However, during actual testing and validation, the model only accesses explanations and text embeddings without the corresponding labels. Both explanations and labels are essential during training. To harness the full capabilities of the pre-trained ExpBERT model, each annotation cycle not only provides labels but also appends new explanations to the default explanation set (Explanation.txt).

Throughout the process, three main sets continuously engage in iterations within three primary modules until the evaluation score based on the **validation set** approaches or exceeds the evaluation score for the full data volume. These modules are training, annotation, and updating:

Training: During the initial training phase, **15%** of the entire data volume, which comes with labels and explanations, is used. Each sample in the ‘exp_label_dataset’ is linked to 9 default explanations (label descriptions) for the first round of training. After the first round of training, the **trained model** and the ‘no_exp_label_dataset’ are used as input parameters for the active learning strategy. The “information content” of each sample is then calculated from perspectives such as uncertainty or diversity. Samples with the highest “information content” are subsequently selected for annotation.

Annotation: In each iteration, the annotator receives **4%** of the samples. During the annotation process depicted in the figure, the annotator identifies key phrases in the text that **semantically resemble** the category description as explanations, and the system automatically provides the classification label for each text. Once samples are fully annotated, they are added to the ‘exp_label_dataset’, forming the dataset required for the next iteration.

Dataset Update: Subsequently, the annotated data, which constitutes **4%** of the total data volume, is removed from the ‘no_exp_label_dataset’. With each iteration, the size of the ‘exp_label_dataset’ increases by **4%**, and one or more explanations are added. After updating the ‘Explanation.txt’, the total number of explanations will exceed the initial default set. All these datasets are pre-trained by the ExpBERT-based model before being fed into the neural network text classifier.

3.2 Pre-training

According to the framework design mentioned above, all datasets require a pre-training process before training the text classifier. As outlined in Section 2.1 about ExpBERT, pre-trained models are typically trained on a large volume of text data, capturing deep bidirectional language representations. As a result, they often achieve higher performance on specific tasks. In many text classification scenarios, using a pre-trained model as a foundation and fine-tuning it can significantly reduce the training time of the text classification model.

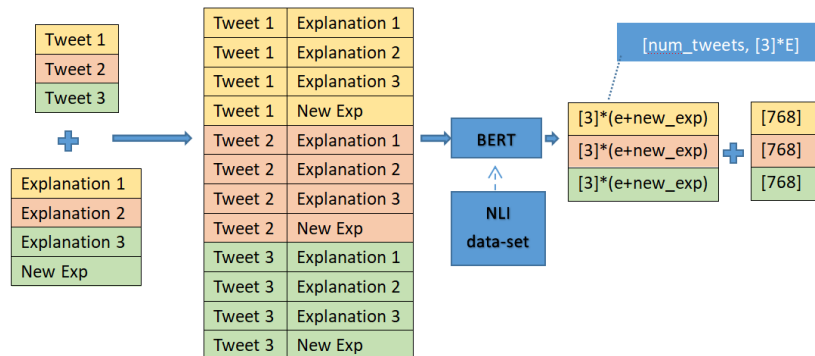


Figure 3.2: Overview of pre-training process

Text is typically input into the model to acquire its embedding representation. This embedding representation can then serve as the input for downstream text classifiers. However, the embedding size

obtained from the original model is often too large, leading to extended computation times. Therefore, as illustrated in Figure 3.2, Natural Language Inference (NLI) is employed to fine-tune BERT, resulting in an embedding size of $3 \times E$ for each sample. To retain tweet information, after concatenating the text embedding representations, the size of a sample’s embedding is only **768 plus $3 \times E$** .

In addition, it is assumed that the number of explanations corresponding to the annotated dataset will increase by ‘new_exp’ during the active learning iterations. Therefore, every time the iteration reaches the pre-training step, it is necessary to re-connect the updated training set, test set, and validation set texts with all the updated explanations as input to the pre-training model. As a result, each pre-training session will handle an additional **num_tweets*new_exp** amount of data. This aspect also contributes to the extended duration of the experiment implementation.

For the unannotated dataset, to simulate real-world conditions, only the text data needs to be processed during pre-training. However, the trained neural network expects to receive input with the same dimensions as during training. Thus, for the embedding representation of unannotated data, the ‘F.pad(embedding, (0, pad_amount))’ method is used to ensure the tensor dimensions of the embedding match. Lastly, pre-training employs a pool-based approach to speed up the computation to enhance the computational efficiency of each iteration. The ‘Pool’ class from Python’s ‘multiprocessing’ library is utilised to process data batches in parallel. Multiple batches can be processed simultaneously on multi-core machines, allowing for flexible process control and increased processing efficiency.

3.3 Classifier Model

This paper employs neural networks as the classifier model, taking the post-pre-training embeddings as input. The evaluation phase (Chapter 5) will compare the performance of conventional neural network classifiers and neural network classifiers with ‘dropout’ in conjunction with active learning techniques. Structural modifications will be made to the traditional neural networks in this experiment. The final model structure will be built using PyTorch [31].

3.3.1 Before Dropout

Regarding model architecture, the neural network in this study employs a single hidden layer and 100 neurons. The output feature count, or ‘output size’, is 9 (number of labels). The forward propagation design can be represented by the following formula, where W_i and b_i are the weights and biases of the linear layer, respectively. After passing through an activation function, the model output y is obtained. This study will use the *ReLU* non-linear activation function, which induces network sparsity, alleviating overfitting issues [71].

$$h = W_i x + b_i, \quad a = \text{ReLU}(h), \quad y = W_{i+1} a + b_{i+1} \quad (3.1)$$

Regarding training methods, the conventional neural network uses the standard loss computation method (cross-entropy [13]). Additionally, the steps of forward propagation, loss computation, backward propagation, and optimisation are executed separately using PyTorch’s standard methods.

3.3.2 After Dropout

This paper has improved the model structure by employing a dropout neural network model as a text classifier when applying the BALD method. Hinton mentioned that Dropout can prevent neurons from overly relying on other specific neurons in the network, thereby enhancing their independence and strengthening the model’s generalisation capability [28]. Applying dropout to neural networks is equivalent to having a “sparse” network made up of units that have survived the dropout. A neural network with n units can have 2^n possible sub-network configurations. For each presented training instance, a new sparse network is sampled and trained [73].

Dropout effectively approximates different neural network structures by randomly discarding units. The dropout process is illustrated in Figure 3.3, achieved by removing some neurons from the network and breaking their connections. In this paper, the ‘dropout_prob’ parameter is set to **0.2** during the

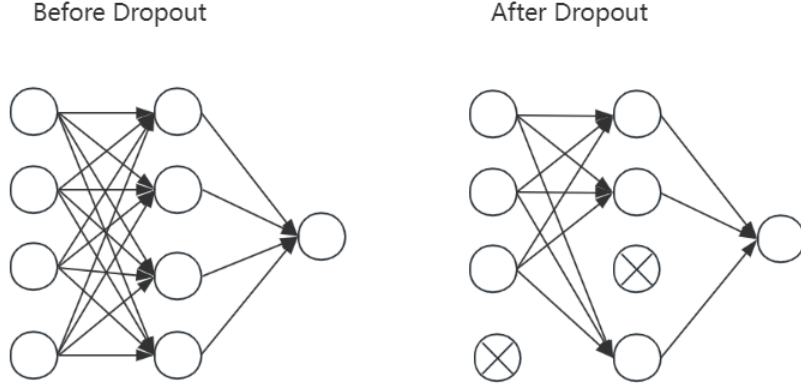


Figure 3.3: Standard Neural Network and Dropout Neural Network.

initialisation of the neural network, meaning that the probability of a neuron being temporarily removed from the network during each forward propagation is 20%. At this point, the forward propagation formula for the original neural network changes:

$$a = \text{ReLU}(h), \quad r = \text{Bernoulli}(p), \quad a_{\text{Dropout}} = r * a, \quad y = W_{i+1} * a_{\text{Dropout}} + b_{i+1} \quad (3.2)$$

$\text{Bernoulli}(p)$ represents a random variable following the **Bernoulli distribution**, used to generate a probability vector r , which randomly yields a vector of 0s and 1s. Consequently, some elements in vector a are set to 0 based on this distribution.

3.4 Active Learning Strategy

In Section 3.1, the acquisition function process was mentioned. In the design, **four different active learning strategies** were adopted to try to improve model performance: random sampling as a baseline; uncertainty sampling based on Least Confidence; diversity sampling based on sentiment; and Bayesian Active Learning by Disagreement (BALD) based on MCD and Least Confidence. Each of these strategies explores and analyses the characteristics of the data in different ways, serving as the core of the active learning framework to enhance the efficiency and accuracy of model learning. The best ways to improve model performance are explored by sampling from different perspectives.

Random Sampling

During those sampling processes, the acquisition function will evaluate nearly 70% of the unannotated dataset, and each evaluation provides 4% of the unannotated dataset to annotators for annotation. Random sampling ensures **unbiased** sample selection. The sampling is done directly on the unannotated dataset using Python's built-in random module without taking the trained model into account. The model performance under its influence will serve as a **baseline**.

Uncertainty Sampling Based on Least Confidence

This paper employs the Least Confidence approach for uncertainty estimation. As discussed in Section 2.3.2, the Least Confidence method is often used as an alternative to entropy computation for multi-label text classification problems. In this experiment, the model is set to evaluation mode. The **highest probability** value is sought within the probability distribution 'prob_dist' given by the model. The number of labels 'num_labels' is calculated, and the most confident prediction is normalised using the formula below. The resulting 'normalized_lc' provides a standardised score where a **higher value indicates less confidence** in the sample.

$$\text{normalized_lc} = (1 - \max(\text{prob_dist})) \times \left(\frac{\text{num_labels}}{\text{num_labels} - 1} \right) \quad (3.3)$$

Samples are ranked in ascending order based on this score, and the top 4% of samples with the **highest scores** are selected as the output for this uncertainty method.

Semantic-based Diversity Sampling

In this study, a semantic diversity-based sampling method is employed. This approach leverages a perspective distinct from uncertainty sampling in the active learning process to identify enriching and minimally redundant samples. For the dataset under consideration, there is a significant imbalance in the number of samples across categories. Some categories account for only 2% of the entire dataset. Such imbalances might lead to certain categories being under-trained, resulting in considerable uncertainty during evaluation. Relying solely on uncertainty for sampling may lead to the over-sampling of a specific category, thereby skewing the distribution of the sampled results, consequently affecting the learning outcome of the model. Hence, a semantic diversity-based sampling strategy is adopted to utilise samples from all categories during active learning effectively.

This method requires loading an unannotated dataset and using the embeddings from a pre-trained model. By setting the trained neural network model to evaluation mode, probabilities corresponding to the input embeddings are computed. As introduced in Section 2.3.2, each text’s embedding vector is denoted as Y_{i*} , where i represents the index of the text. Suppose the dataset is divided into n batches, with m instances each. Initially, a random subset of embedding vectors is chosen as the clustering centres, denoted as C_{i*} . In this study, C_{i*} is represented as follows, where the proportion of the initial clustering centres, or α , is set to 0.2:

$$C_{i*} = \text{random_choice}(Y_{i*}, \alpha * n * m) \quad (3.4)$$

Subsequently, the greedy k-centre algorithm is utilised to search for embedding vectors in Y_i that aren’t in C_{i*} . The algorithm identifies embedding vectors that are the **furthest from the existing clustering centres** and updates the set of clustering centres. This process is illustrated as:

$$C_{i*} = \text{argmax}_{y \in Y_{i*} \setminus C} \min_{c \in C_{i*}} \text{dist}(y, c) \quad (3.5)$$

The ‘dist’ function computes the **semantic distance** between embedding vectors and clustering centres. Finally, we return the indices in Y_{i*} for each cluster centre in C_{i*} denoted as O_i .

MCD Least Confidence BALD

This study adopts a sampling strategy based on Monte Carlo Dropout(MCD) to find samples with higher model uncertainty for more effective model training. The foundation of this strategy is a study by Gal and Ghahramani in 2016 [20], where they proposed using dropout during neural network training as an approximation to Bayesian inference. Similar to other active learning, the embedding representation of the unlabeled dataset is loaded.

Algorithm 3.1: BALD using MCD

```

1 Set model to TRAINING mode
2 for input_data in Embeddings do
3   Initialise predictions as an empty list
4   for  $j = 1$  upto  $T\_samples$  do
5     logits  $\leftarrow$  model.forward(input_data)
6     probabilities  $\leftarrow$  SOFTMAX(logits)
7     APPEND probabilities to predictions
8   end
9   avg_probabilities  $\leftarrow$  AVERAGE(predictions)
10  bald_score  $\leftarrow$  NEGATIVE MAXIMUM of avg_probabilities
11  APPEND (input_data’s index, bald_score) to selected_indices
12 end
13 SORT selected_indices by bald_score in DESCENDING order
14 return top k indices from selected_indices

```

However, the difference is that the text classification model is a neural network with dropout. Also, as shown in Algorithm 3.1, it needs to be **set in training mode**, not evaluation mode. This is because MCD requires the dropout to be active during training. Unlike the traditional BALD strategy, the experiment measures uncertainty by **maximising the negative maximum of the prediction probability** rather than the conventional BALD calculation method. Next, the indices of all samples and their corresponding

BALD scores are stored in a list and sorted in descending order based on the BALD scores. In this way, we can identify the samples for which the model is most uncertain about their classification, i.e., samples with the **highest** BALD scores. This method can be seen as a variant of the BALD algorithm that integrates MCD and the Least Confidence strategy. Compared to uncertainty sampling that only uses the Least Confidence, this method can capture **more stable uncertainty information** from the model.

3.5 Annotation Process

After the sampling guided by the active learning strategy, a set of unannotated texts with high information content is obtained. For these texts, annotators need to provide the corresponding classification based on their content and decide which explanatory notes to add in the current iteration.

3.5.1 Label Annotation

The focus of this experiment is to explore whether the addition of explanatory annotations during the active learning process can enhance the performance of the text classifier. To simulate an ideal situation, it is assumed that in each iteration, the annotators have a complete understanding of the text classification, so the provided labels will **accurately correspond** to the source data. Furthermore, in the evaluation stage of Chapter 5, to eliminate the potential impact of the added labels and text volume on model performance, we will compare two scenarios: one with active learning that only added labels and another with active learning that added both labels and explanations.

3.5.2 Explanation Annotation

This experiment designed three different explanation annotation strategies. The provided explanations will be stored in the ‘Explanation.txt’ set displayed in Figure 3.1, available for the text concatenation task prior to subsequent pre-training. System operators can choose among these three methods based on workforce allocation or cost control to better support the model’s learning process:

The first strategy involves previewing the texts in each iteration in advance and simulating the possible added explanations, just as an annotator would. The **preset explanations** will be saved in a ‘txt’ file. When the annotation phase is executed, explanations are automatically extracted from this file. Although this method can reduce workforce consumption, it lacks adaptability and is only suitable for the dataset related to this experiment.

The second strategy allows **humans** to provide input as explanations directly. This approach requires users to analyse the features of the selected text and provide their explanations. In each iteration, the annotator will focus on ten prominent texts of a particular category to quickly analyse these texts. This method incorporates human judgment and applies to other emergency-related data, offering some flexibility.

The last strategy uses **OpenAI’s model** to automatically provide explanations, reducing the need for human effort. This experiment chose the ‘**text-davinci-003**’ engine for text analysis [78]. Through the developer platform’s **Playground** on the OpenAI interface, the prompt to be input into the model can be simulated in the Playground. The prompt uses “*generate a string that only contains the words in the text that are only semantically equal to*” to connect the active learning strategy-selected prominent texts of a particular category. On the right navigation bar of the Playground, one can select the ‘Model’, ‘Temperature’, and maximum response length. The best-performing hyperparameters are chosen based on the reasonableness of the response. As shown in Table 3.1, for the classification “deaths and injured”, tests found that the ‘text-davinci-003’ responses were far more reasonable and accurate than those of “text-ada-001”. Other configurations tended to generate irrelevant information like geographical locations and numbers.

After debugging, it was found that when the ‘**Temperature**’ (the lower the temperature, the less randomness) is set to 1 and ‘**Top_P**’ (which controls diversity via nucleus sampling) is set to 1, the

Model	Temperature	Top P	Response
text-ada-001	0.5	1	“As many as 98% of those killed .. <i>Pakistan Somalia</i> ”
text-ada-001	1.0	1	“Hundreds dead in <i>Pakistan</i> .. country 6.0 quake, dozens hurt”
text-ada-001	1.0	0.5	“As many as 98% of those killed strikes are civilians.”
text-ada-001	1.5	1	“patrol western shuts the exchange <i>peace 65 party</i> ”
text-ada-003	0.5	1	“Hurt Death <i>Pakistan</i> Quake Dozens Killed Loved Ones”
text-ada-003	1.0	1	“dozens hurt killed families loved ones RIP injured”
text-ada-003	1.0	0.5	“Hurt, Dozens, Killed, <i>Pakistan</i> , Families, Loved, Injured”
text-ada-003	1.5	1	“death Mortality over 320 <i>people</i> killed drone 98% victims”

Table 3.1: Results of the OpenAI model evaluation.

model’s output is most reasonable. The Playground provides configuration code for the Python environment, allowing the code to be directly applied to the system. To successfully connect to the OpenAI interface, it is necessary to apply for an **OpenAI key** and load this key before generating responses.

It can be observed that using this model can analyse texts, selecting keywords that are semantically closest to the category label description to concatenate into a string. Ultimately, it returns a string that best reflects the characteristics of the text and helps the model learn from it as an added explanation. However, its drawback is that it’s not friendly for the demand of generating multiple explanations at once or for active learning of hundreds of tokens. Although it reduces workforce costs, the generation of each token consumes certain funds.

Chapter 4

Implementation

This chapter primarily discusses the setup required to conduct the experiment. Section 4.1 introduces the fundamental environment in which the experiment is conducted, and Section 4.2 details the preparation and splitting of the dataset used in the experiment. Additionally, the performance metrics (Section 4.3) used in the evaluation in Chapter 5 and the rationale behind their selection are explained. Lastly, the stopping criteria adopted for the experiment (Section 4.4) and the execution modules (Section 4.5) are specified.

4.1 Experimental Environment

This experiment employs **Python 3.9** as the primary development language for natural language processing tasks. Being highly suitable for natural language processing tasks, Python boasts a plethora of data science and machine learning libraries [33]. This study utilised the capabilities provided by libraries such as PyTorch, Pandas [46], and Numpy [55] for data analysis, pre-processing, and machine learning. Python’s high compatibility facilitates the connection to remote High-Performance Computing Systems (HPC) within the context of this study using PyCharm, or running through Google Colab in conjunction with Google Drive. This research integrates **PyTorch** with Python to construct dynamic neural networks. By leveraging PyTorch’s automatic gradient computation, the features of Dynamic Computational Graph, and the capability of CUDA for accelerated computation, a high level of flexibility and accuracy is maintained in the neural network model [57].

4.2 Dataset

The dataset used in this experiment simulates the text distribution of social media in emergencies. The dataset is prepared and split according to the framework described in section 3.1 to meet the training requirements with a small amount of data in active learning.

4.2.1 CrisisNLP

In this experiment, the CrisisNLP dataset is utilised, which comprises data IDs, textual information, and descriptions of their corresponding categories. There are nearly 16,000 entries in total. This dataset requires preprocessing, splitting, pre-training, and classification. CrisisNLP is an open-source resource designed explicitly for humanitarian emergency and crisis response related natural language processing (NLP) research and development [32]. The dataset contains interaction information from various large social media platforms, covering multiple crisis events, such as typhoons, earthquakes, and other events that require timely responses from relevant departments. The nature of the information can be divided into nine categories. The descriptions of these labels and their distribution are shown in the Table 4.1.

According to the distribution in Table 4.1, it is evident that the data proportion is **highly imbalanced**. Label 7 has the highest proportion, and the information it provides is also the most mixed,

Label	Description	Percentage
0	injured_or_dead_people	13%
1	missing_trapped_or_found_people	2%
2	displaced_people_and_evacuations	3%
3	infrastructure_and_utilities_damage	8%
4	donation_needs_or_offers_or_volunteering_services	14%
5	caution_and_advice	6%
6	sympathy_and_emotional_support	11%
7	other_useful_information	30%
8	not_related_or_irrelevant	13%

Table 4.1: Data distribution in CrisisNLP.

accurately reflecting the distribution of emergency information in real-life scenarios. This poses a challenge for text classification tasks. Therefore, this dataset is of great value to researchers, as it documents real-time reactions and behaviours of people in various crisis situations and their corresponding data distribution. Consequently, researchers can develop more effective systems tailored for genuine data imbalances, addressing these crises, such as automatically detecting crisis-related social media posts or predicting the development trend of a crisis through analysing social media data.

4.2.2 Data Preparation

Preprocessing: Before feeding into the pre-trained model, preprocessing the tweets in the dataset is a critical step. The goal is to minimise the noise in the text and enhance the model’s learning efficiency. Firstly, internet slang is standardised, converting informal abbreviations or colloquialisms into their standard forms and splitting camel-case named words. Subsequently, emojis in tweets are replaced with their textual counterparts using the emoji library, enabling the model to better understand these symbols’ meanings. Lastly, website links, usernames, and “#” tags in tweets are cleaned or replaced since these elements typically don’t contribute much to the model’s learning. Such preprocessing steps help in reducing noise in tweets, thereby enhancing the model’s learning efficiency and performance.

Dataset Splitting: For active learning tasks, splitting the dataset into unannotated and annotated datasets is essential. At the same time, the annotated dataset will be further split into training, validation, and testing sets. Active learning requires adding or removing data in each iteration. However, to ensure that the model’s performance on unseen data can be tested fairly in each active learning iteration, the number of test and validation samples must **remain constant**. Thus, the test and validation set, unannotated and annotated datasets must be split **outside of the active learning loop**. The **training set will expand** with each iteration as the dataset grows. Since the data corresponding to label 1 only account for about 2% of the total dataset, to ensure that all data subsets contain data corresponding to label 1, we use the ‘StratifiedShuffleSplit’ mechanism from the scikit-learn library to maintain the original data’s class distribution. Compared to random splitting, this method provides a more accurate assessment of the classifier’s performance.

Data	Count	Percentage
Annotated Data	Training dataset	2408
	Validation dataset	803
	Test dataset	1606
Unannotated Data	11241	70%

Table 4.2: Distribution of datasets

To meet the criteria of active learning, which achieves high performance with a small amount of training data, the number of annotated data in the first iteration of active learning is not the total data. As shown in Table 4.2 above, the total is **only 30%** of the entire dataset, with a length of 4827. Therefore, the size of the training set is no longer the most significant component, with the training set representing only **15%** of the total data for the active learning strategy and the size is set to 2408. The size of the unannotated dataset is set to **11241**, which accounts for 70% of the total dataset. The unannotated

dataset becomes the most significant component compared with other datasets.

When training with the **full dataset** (with the training set having the highest proportion), the combined length of the training and validation datasets in the annotated data will be 12042. The number of the test set remains unchanged both in the full dataset and in the amount of data required for active learning, with a length of 1606. 4% of the data is sampled for annotation in each cycle. In the subsequent evaluation chapter, it was found that after 9 iterations, the model performance approached the model’s performance corresponding to the validation set under the full data. At this time, the length of the unannotated dataset remains 6741.

4.3 Performance Metrics

The setting of performance metrics can significantly influence the effectiveness of active learning projects. Proper evaluation criteria can measure and compare the impacts of different active learning strategies and model frameworks on emergency systems in human learning. This, in turn, determines the choices of models, active learning strategies, and parameters. To assess the performance of active learning, the data is first partitioned according to the standards in Section 4.2.2. Then, the training begins on the training set. After each pool-based active learning iteration round, the trained model is evaluated on the test and validation sets using the predefined performance metrics. After each iteration, learning curves are plotted based on the performance metrics.

To comprehensively assess the model’s performance, the experiment utilises two performance metrics: **accuracy and F1 score**. Accuracy measures the proportion of instances that the model correctly predicts relative to the total number of instances, and it can be a biased evaluation criterion [18]. In contrast, the F1 score is the harmonic mean of precision and recall. Compared to accuracy, it eliminates bias and is more commonly used in multi-class problems [26].

$$F1_score = \frac{2}{(1/Precision + 1/Recall)} \quad (4.1)$$

Precision is the proportion of positive instances correctly predicted by the model out of all instances predicted as positive, and Recall is the proportion of actual positive instances that the model correctly predicts as positive. When dealing with imbalanced datasets like CrisisNLP, the F1 score can provide a more comprehensive measure of performance.

In addition, the **learning curve** can also serve as a performance evaluation criterion. By observing the changes in the learning curve, we can better understand the model’s learning capacity and stability [24]. Tensorboard [56] records the learning curve during the active learning iteration process. The learning curve depicts the relationship between the model’s performance and the number of samples used for training (number of active learning iterations), evaluating whether increasing the number of annotated datasets can effectively enhance the model’s performance. Moreover, the experiment also logs the performance of the neural network classifier for each epoch during each training process through Tensorboard, ensuring the model achieves the best performance without overfitting after each active learning iteration.

4.4 Stopping Strategy

Within the active learning framework, the Stopping strategy is used to determine when to halt the learning process or no longer request new data from annotators [39]. The design of the Stopping strategy directly affects the cost of annotation, time consumption, and the final performance of the model. Although the Stopping strategy is somewhat akin to a model’s hyperparameters, it doesn’t directly influence the model’s training process. Instead, its primary objective is to reduce annotation costs and save time while maintaining optimised performance. This strategy aims to balance time efficiency and performance output to achieve a learning process that is both efficient and high-performing.

For the training of neural network models, the configuration of hyperparameters directly impacts the training process and model performance. In early research, there have been studies that meticulously

explored and determined various hyperparameters. Here, we referenced these studies and selected hyperparameters suitable for neural network training, namely using a batch size of 8 and a learning rate of $5e-5$.

When devising the stopping strategy for the active learning experiment, monitoring whether the model performance has reached saturation is essential. For this experiment, the primary focus is on the average F1 score of the model on the **validation set**. Specifically, we aim to observe when the model, based on a random sampling strategy, can surpass the performance of the model trained on the entire dataset or when the rate of performance growth begins to decelerate.

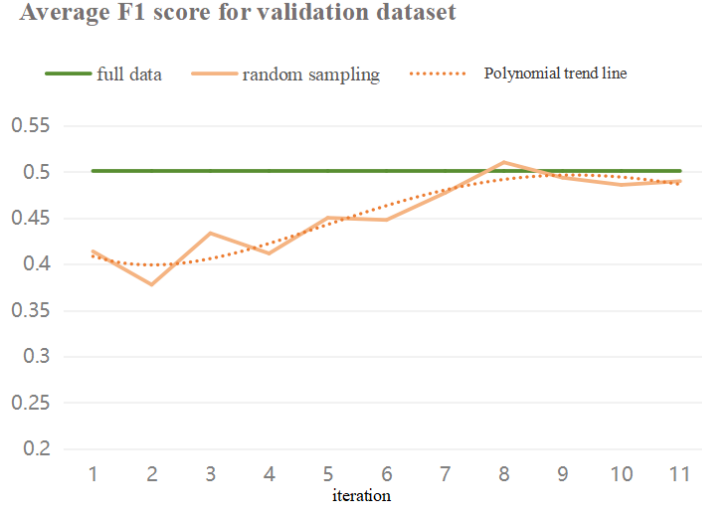


Figure 4.1: Average F1 score of full-data and random sampling strategies under validation dataset

Figure 4.1 reflects the performance comparison between the models trained under the **random sampling** and **full-data** training strategies. By the 8th iteration, the model’s performance based on random sampling had already exceeded that of the model trained on the full data. However, this performance improvement was inconsistent, so we continued iterating. After the **9th iteration**, the growth in model performance began to slow down. Considering that further iterations would increase time and annotation costs, we terminated the experiment after the 9th iteration. This strategy allows for **achieving optimal performance at minimal costs**.

4.5 Execution Module

In this study, we conducted text classification by implementing a series of active learning strategies. The experiment is run either by connecting to a remote environment through PyCharm or locally. The operational source code, as well as the data, is stored in a version control tool. The address of the repository is shown in Appendix A. In most cases, an explanation is added in each iteration, with a total of 9 iterations being performed. At the conclusion of these iterations, the experiment for each active learning strategy took a total of **24 hours**. Using the MCD combined with the BALD active learning strategy, the time taken approaches 37 hours. The experimental execution module is primarily segmented based on active learning strategies. As the evaluation in Chapter 5 found that diversity sampling performed the best, different explanation annotation methods were employed to further explore the impact of explanations on performance.

- Random Sampling: Implemented through the “Whole_train.py” script, processing both full data mode and active learning mode.
- Senmatic-based Diversity Sampling: Based on the “diversity_sampling-whole.py” script, this method employs three different explanation annotation methods: preset explanations, manual input, and explanations generated by the OpenAI model (**an API key for OpenAI is required**).

- **Uncertainty Sampling:** Using the “uncertaintyWhole.py” script, the aim of this strategy is to select data points that present the most uncertainty or challenge for model classification by calculating the minimum confidence level.
- **MCD-based BALD Strategy:** Executed via the “BALD_MCD.py” script, it combines dropout neural networks with the Bayesian Active Learning by Disagreement (BALD) strategy, specifically seeking high uncertainty data points in model predictions through minimum confidence.

The primary focus of the experimental execution is to conduct a deep evaluation of the performance of various strategies in text classification tasks. This aims to ascertain which strategy and explanation annotation method is best suited to enhance the model’s learning efficiency.

Chapter 5

Critical Evaluation

This Chapter first sets out the expectations and framework for evaluation (Section 5.1). Based on five expectations, combined with the evaluation framework and the performance metrics provided in Section 4.3, the evaluation is carried out (Sections 5.2, 5.3, 5.4, 5.5, and 5.6). Results corresponding to each expectation are presented, and the evaluation results are analysed to see if the outcomes align with the preset expectations. Lastly, Section 5.7 provides a summary of the evaluation results.

5.1 Expectations and Evaluation Framework

Initially, there are several preset expectations for the Human-in-the-loop framework.

- Pool-based active learning can effectively enhance model performance (on both test and validation sets) by increasing the size of the annotated dataset.
- The addition of explanation annotations during the annotating process has a greater impact than using data with label-only annotations.
- Most strategies can surpass the baseline (random sampling) while exceeding or approaching the full-data dataset’s performance before implementing the stopping strategy.
- We can identify the best active learning strategy, explanation annotation strategy, and model architectural strategy to quickly achieve high performance with a minimal amount of training data.
- The number and quality of the explanations provided greatly influence model performance.

Notably, the full-data dataset refers to the annotated dataset (with nine default explanations linked), where the training set makes up 72% of the data distribution, and the remaining percentage is divided between the validation and test sets. However, the training set share in active learning is 15%.

Based on these expectations, an evaluation framework is established in Table 5.1, which will be applied to both the test and validation sets using F1 score and accuracy as the performance metrics:

ID	Model structure	Dataset	AL strategy	annotation method
1	No dropout	full-data dataset	None	without explanation
2	No dropout	AL-based datasets	Random sampling	without explanation
3	No dropout	AL-based datasets	Random sampling	explanation-NoOpenAI
4	No dropout	AL-based datasets	Uncertainty sampling	explanation-NoOpenAI
5	No dropout	AL-based datasets	Diversity sampling	explanation-NoOpenAI
6	No dropout	AL-based datasets	Diversity sampling	explanation-WithOpenAI
7	With dropout	AL-based datasets	BALD	explanation-NoOpenAI
8	No dropout	AL-based datasets	Diversity sampling	two explanations each iteration

Table 5.1: Evaluation structure.

5.2 Expectation One

According to the expectations set in Section 5.1, a point-by-point evaluation is carried out. For the first expectation, it's necessary to observe whether active learning, by increasing the annotated dataset, can effectively enhance the model's performance. Hence, the learning curve mentioned in Section 4.3 needs to be inspected. The expectation is that within the learning curve, the model's performance consistently rises as the number of annotated instances increases. To reduce computational time, a representative evaluation method was chosen, specifically the one with **evaluation ID 2** from the evaluation framework (Table 5.1) in Section 5.1, employing a random sampling strategy. In each round, only the corresponding classification label is annotated without adding new explanations, instead using the default nine explanations.

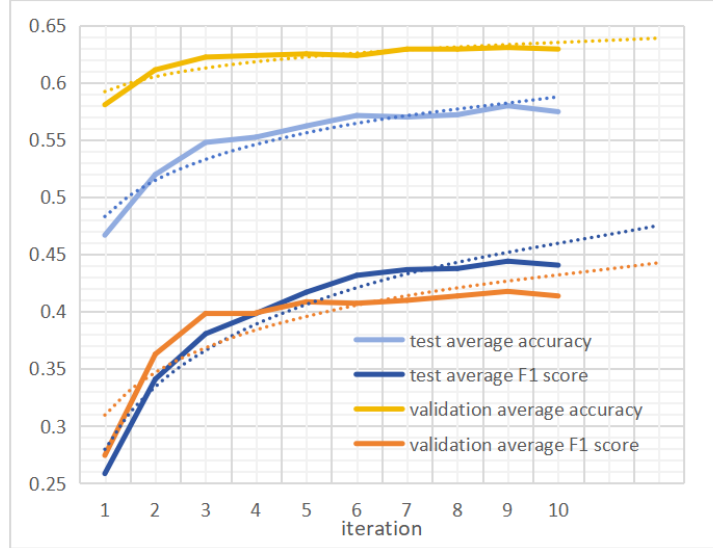


Figure 5.1: Random sampling only annotates labels

The evaluation results are illustrated in Figure 5.1. The lines above represent the average accuracy for the validation and test sets, while the lines below represent the average F1 score. By setting a logarithmic trendline, it's evident that the learning curve's trajectory **aligns with the first expectation**. Simultaneously, the data's accuracy on the validation set is generally **better** than on the test set.

This is due to active learning, as mentioned in Section 2.3, continuously supplying a portion of the annotated dataset to the original dataset for training. As a result, with the increase in input data volume, the classifier can learn more effectively from the abundant data and enhance its generalisation ability. The reason the accuracy on the validation set is better than that on the test set is evident from Table 4.2, which shows that the overall quantity of the validation set is smaller. Hence, there is less noise and fewer outliers. To simulate real-world scenarios, the test set is larger and noisier, leading to slightly lower overall accuracy. Besides, the test set covers a larger distribution of data, and the F1 score could still be higher due to better performance in a few classes.

In summary, this Human-in-the-loop pool-based active learning framework indeed effectively boosts the model's performance, whether on the test set or the validation set.

5.3 Expectation Two

In Section 5.1, the second expectation investigates the effectiveness of the Human-in-the-loop framework when based on ExpBERT. It explores whether, during the annotation process, adding explanations in addition to labels can lead to higher performance. Thus, we utilise the evaluation methods corresponding to **evaluation IDs 2 and 3** from Table 5.1 for a comparative experiment. We examine whether the results of evaluation ID 3 (**with added explanations**) are superior to those of evaluation ID 2 (**with**

only added labels). In the evaluation process for ID 3, the number of explanations increased by one in each round. Here, the evaluation results use the average F1 scores and average accuracy scores for the validation and test sets. The performance from the final iteration is displayed in Table 5.2.

	Average F1 score	Average accuracy score
Evaluation 2 for validation set	0.4136	0.6292
Evaluation 3 for validation set	<u>0.625</u>	<u>0.6131</u>
Evaluation 2 for test set	0.4405	0.5745
Evaluation 3 for test set	<u>0.5056</u>	<u>0.6216</u>

Table 5.2: Results of evaluation ID 2 and evaluation ID 3

From Table 5.2, it is evident that all evaluation results for ID 3 (where explanations were added along with label annotations) are **higher** than those for ID 2 (where only labels were annotated). This is because, as mentioned in Section 1.1, **explanations** include keywords and sentences related to classification information. Adding these explanation links allows the model to learn from a deeper semantic level. After pre-training with the ExpBERT mechanism, the received embeddings enhance the classifier’s generalisation ability. Therefore, the model’s performance is higher in the evaluation framework for ID 3.

Thus, annotating with not just labels but also adding explanation annotations can effectively enhance the classification model’s performance. The conclusion is that the second expectation has been met.

5.4 Expectation Three

Expectation three primarily examines the effectiveness of active learning strategies, as well as whether they can surpass the performance of the baseline (random sampling) and whether they can approach or exceed the performance under the full data set using fewer data before the termination strategy. Therefore, the performance of different strategies during the last iteration of active learning was adopted, i.e., using the evaluation methods corresponding to **IDs 1 (Full data), 3 (Random sampling), 4 (Uncertainty sampling), 5 (Diversity Sampling), and 7 (BALD)** in Table 5.1 from Section 5.1. The baseline, all active learning strategies, and full data performance were compared. Both the average F1 score and the average accuracy rate were used for evaluation. The evaluation results are shown in Table 5.3 below.

	ID1 (full data)	ID3 (baseline)	ID4	ID5	ID7
test_avg_acc	0.6052	0.6216	0.6256	0.633	0.6215
test_avg_F1	0.4827	0.5056	0.51	0.5466	0.5329
val_avg_acc	0.6012	0.6131	0.6026	0.6348	0.6275
val_avg_F1	0.4866	0.486	0.4981	0.5522	0.5286

Table 5.3: Performance of different strategies through iterations

Results: From Table 5.3, it is evident that uncertainty sampling (ID4) and diversity sampling (ID5) outperform the baseline in terms of test set performance. In the validation set, the performance post the implementation of diversity sampling and the BALD strategy surpasses the baseline. Notably, diversity sampling exceeds the baseline across different datasets and is based on various performance metrics. Moreover, the average accuracy of the test sets for uncertainty sampling and diversity sampling is higher than that of the full-data dataset. What’s more, the three metrics related to diversity sampling all surpass the performance metrics of the full dataset. The average F1 score in the validation set under the BALD strategy also outperforms the full-data dataset’s score.

Reason: The reason for such results is that other strategies are biased compared to the baseline [10]. Uncertainty sampling assists the model in learning from the most uncertain data, thereby overcoming

the model’s weaknesses in such data. Conversely, diversity sampling opts for semantically diverse data, allowing for a more comprehensive model training. BALD identifies the data with the highest uncertainty through multiple samplings, enhancing performance. Therefore, compared to unbiased random sampling, the other three active learning strategies can aid the model’s training from various biased perspectives.

In Conclusion, given that most active learning strategies surpass the baseline, and some even exceed the performance metrics of the entire dataset, it can be inferred that Expectation three is met. Consequently, adopting advanced active learning strategies can further improve model performance. It’s even possible to surpass the performance of the full-data dataset with training on just a fraction of the data.

5.5 Expectation Four

Expectation Four primarily aims to select the most suitable solution for the emergency response system to enhance its performance. It requires the selection and comparison of different sampling strategies and analysing of the results produced by these strategies to determine the causes for such outcomes, ensuring they are not coincidental. Firstly, to rule out randomness, learning curves corresponding to these active learning strategies will be plotted to observe the performance at every stage. The test metrics are the same as those in section 5.4. However, instead of just observing the final iteration, we will be looking into every iteration.

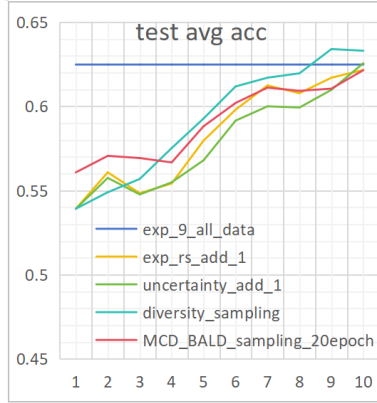


Figure 5.2: Average accuracy score for test set

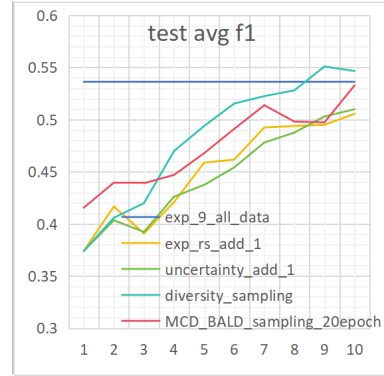


Figure 5.3: Average F1 score for test set

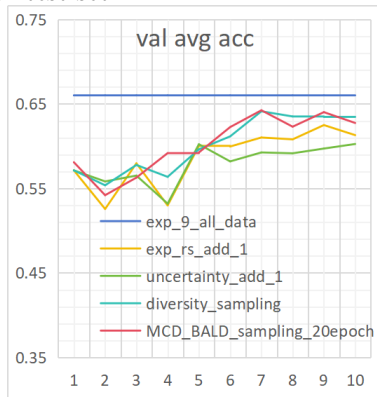


Figure 5.4: Average accuracy score for validation set

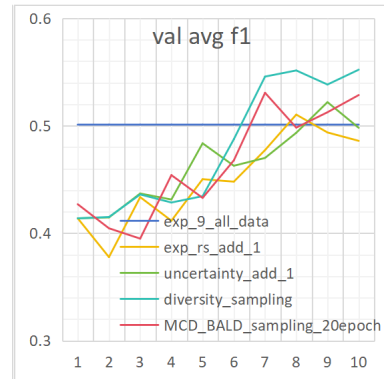


Figure 5.5: Average F1 score for validation set

Results: Four line charts (Figure 5.2, Figure 5.3, Figure 5.4 and Figure 5.5) have been plotted, as shown below, representing the average F1 scores and average accuracy on the test and validation sets under different strategies. From the line charts, the top line represents the performance of the full-data dataset. The learning curves corresponding to diversity sampling and the BALD strategy are observed

to approach or even surpass this line with each iteration. In contrast, the uncertainty sampling strategy scores below the baseline (random sampling) for the first 9 and 10 iterations across the four metrics. However, by the final iteration, the uncertainty sampling strategy surpassed the baseline on the test set.

Furthermore, it is essential to observe the performance of these active learning strategies corresponding to different labels, identifying the most suitable learning strategy for the majority of the labels. The evaluation was done using the F1 scores corresponding to the results of the last iteration on the test set. The final evaluation results are shown in Figure 5.4 below. The diversity sampling method yields higher F1 scores for most of the categories, particularly for categories 2, 3, 4, and 7. Uncertainty sampling scores higher for label 5, while the BALD strategy performs better for category 1.

Class label	ID3 (baseline)	ID4	ID5	ID7
0	0.6248	0.6216	0.6256	0.633
1	0.6248	0.6216	0.6256	0.633
2	0.6248	0.6216	0.6256	0.633
3	0.6248	0.6216	0.6256	0.633
4	0.6248	0.6216	0.6256	0.633
5	0.6248	0.6216	0.6256	0.633
6	0.6248	0.6216	0.6256	0.633
7	0.6248	0.6216	0.6256	0.633
8	0.6248	0.6216	0.6256	0.633

Table 5.4: Performance of different strategies in the last iteration

Reason: As mentioned in Section 4.2, the distribution of labels in the dataset is extremely imbalanced. Therefore, in scenarios where data distribution is highly imbalanced, the semantic-based diversity sampling strategy performs best. The underlying reason is that this strategy selects samples that cover a variety of categories with diverse and representative information, providing the model with a comprehensive perspective [16]. This is especially important for imbalanced data because some less frequent categories might be overlooked or undersampled.

Furthermore, combining the BALD and least confidence strategies can also achieve good results. This is because BALD focuses on selecting data with the highest uncertainty from multiple samplings, working in tandem with the least confidence strategy to identify the data points where the model is most uncertain. However, purely using the least confidence uncertainty sampling might not yield the model’s optimal performance in this scenario. This algorithm primarily focuses on the data points where the model is most uncertain, but in highly imbalanced data, this could lead the model to continuously sample from the dominant categories, neglecting the rarer ones.

In summary, considering the overall performance, semantic-based diversity sampling is more suitable for improving the model’s learning rate and overall performance within the active learning cycle in this experiment. Using semantic-based uncertainty sampling can also easily surpass the full data performance after a few iterations.

5.6 Expectation Five

In Section 5.4, we concluded that semantic-based diversity sampling can most effectively enhance the model’s learning rate. Therefore, we further explored the impact of changing the quality and quantity of explanations during the annotation process on overall performance based on this strategy. We selected the evaluations with ID 5 and 6 from Section 5.1 for quality comparison. ID 6, like ID 5, also employed the diversity sampling method. The difference is that in the evaluation corresponding to ID 6, explanations were generated using the OpenAI model mentioned in Section 3.5. It is worth noting that the explanations generated by the OpenAI model are nearly rational, but compared to the preset explanations in ID 5, there are certain imperfections. Due to the word limit of the ‘prompt’ when generating explanations with the OpenAI model, the model couldn’t thoroughly analyze all the texts. Therefore, the

generated explanations often contain a lot of noise, such as irrelevant words describing places, characters, etc., unrelated to the classification description.

Impact on Quality: Hence, the process of generating explanations using the OpenAI model can be regarded as a low-quality explanation annotation process. Therefore, by comparing the differences between these two evaluation methods, we can observe the impact of explanation quality. As shown in the following Table 5.5, during the last iteration of active learning, the evaluation method with ID 5 had performance metrics approximately 2% higher across all test sets compared to the evaluation method with ID 6. Even though the model displayed commendable robustness (with a small performance gap), the quality of the explanations still caused a disparity in performance.

	Average F1 score	Average accuracy
ID5 (validation)	0.6248	0.6216
ID6 (validation)	0.6248	0.6216
ID5 (test)	0.6248	0.6216
ID6 (test)	0.6248	0.6216

Table 5.5: Results of evaluation ID 5 and evaluation ID 6

Thus, using high-quality explanations during the annotator’s annotation stage of the active learning cycle can effectively improve performance. However, considering the potential human resource costs of high-quality explanations, if the aim is to save on these costs, and a performance difference of around 2% compared to high-quality explanations is acceptable, then the approach of generating explanations using OpenAI might be worth considering.

Impact on Quantity: Suppose the amount of new explanations used for annotation in each round increases to 2, and 4% of the data is still provided for annotation in each round. Therefore, the expected effect is to surpass the baseline and achieve the performance under full data with fewer iterations. However, the cost and time required for the ExpBERT model to train in each round also increase accordingly. It is also essential to observe whether the time consumed to achieve the desired performance is equal to or less than the time taken in situations where the number of new explanations added in each iteration is 1.

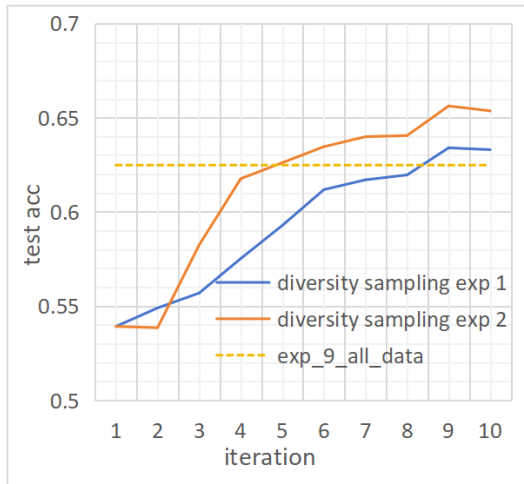


Figure 5.6: Average accuracy score in test set

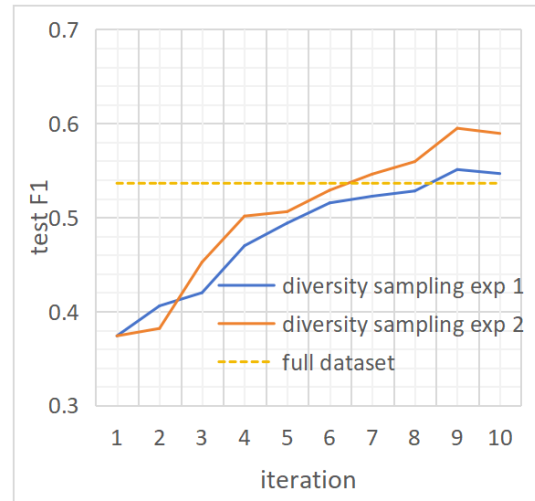


Figure 5.7: Average F1 score in test set

From the image on the left, it can be observed that, by the sixth iteration, the accuracy of the test set with added explanations has already reached or even surpassed the accuracy of the test set with the original annotations. At this point, the runtime of the active learning framework is only around 18 hours. If the emergency response system requires only high accuracy, the termination strategy can be adjusted according to the need. If it continues to run, the **final accuracy can reach 0.65**. Similarly, the right

side indicates that the F1 score of the test set also shows a performance catch-up in the sixth and seventh iterations. However, according to the termination strategy, interrupting the cycle when performance surpasses can reduce the response time of the emergency system.

Therefore, adding an explanation with each iteration can shorten the system’s response time and improve accuracy. Since this method uses fewer iterations (based on the termination strategy, it can be terminated after the sixth iteration), the added data is also reduced, thereby saving on the annotator’s efforts. This approach cleverly balances runtime, workforce, and performance.

5.7 Conclusion

Upon evaluation, active learning can effectively enhance the performance of models and reduce workforce and annotation costs. At the same time, adding explanations in addition to labels during the annotation process has significantly improved the post-annotation model performance. By combining various advanced active learning strategies, it’s possible to effectively surpass the baseline (random sampling) and even outperform the full data performance with fewer iterations and a smaller training set. The most outstanding performance method is to adopt semantic-based diversity sampling combined with high-quality explanatory support. Under the premise of controlling annotation costs and pre-training time consumption, slightly increasing the number of explanations per iteration can achieve the model’s optimal performance.

Chapter 6

Conclusion

This thesis explores the efficacy of a text classifier based on the ExpBERT model applied to disaster emergencies within a Human-in-the-loop (HITL) framework. In simple terms, it seeks the optimal solution that enhances performance while reducing manual annotation costs and computational time. The project encompasses the entire lifecycle: initially designing the framework and the complete active learning process, then constructing a highly compatible experimental environment and providing stopping strategies and data, and finally, evaluating the key modules that influence the system and validating the effectiveness of active learning and explanation annotations. It assesses the impact of different active learning strategies and ways of explanation annotation, as well as balances the quality and quantity of explanations. In this Chapter, we will validate and summarise the entire project. Additionally, this Chapter will address the limitations introduced by this framework and the work needed in the future.

6.1 Project Status

The primary framework of the project integrates humans in the loop with text classification functionalities to enhance the performance of text classification tasks. Different active learning strategies, annotation methods, and model structures are employed to improve model performance. The eventual implementation of the project achieved all the objectives outlined in section 1.3.1. These objectives led to several outcomes. The goals are as follows:

- Develop a pool-based Human-in-the-loop framework applied to the ExpBERT-based text classifier and investigate its effectiveness.
- Construct an annotator simulation process and design stopping criteria (e.g., terminate when reaching the performance of the full data model).
- Establish a baseline (performance of the full-data model without active learning or random sampling) and study whether active learning can enhance the performance of the text classification model. Compare the effects of different acquisition functions.
- Explore the impact of the quality and quantity of explanations on this system. For instance, we provide explanations with noise or increase the number of explanations added in each iteration.

The framework for the first objective was elaborately designed in Chapter 3, and its effectiveness was validated in Sections 5.2 and 5.3. The pool-based active learning framework, both on the validation and test sets, continually showed an upward trend in the learning curve. This Human-in-the-loop active learning framework effectively enhanced the model's performance.

The second objective was realised through the three different explanation annotation methods in Section 3.5 and the stopping parameter settings in Section 4.4. Among the various annotation explanation methods, Section 5.6 found that pre-set explanations or manually annotated high-quality explanations led to a performance 2% higher than those generated by the OpenAI model. This slight difference of 2% also indicates that the OpenAI model-generated explanations can be an alternative when the cost of

manual annotation is substantial.

The realisation of the third objective has a significant impact on performance. Section 3.4 designed the performance lower limit and various biased active learning strategies. It was discovered that models using semantic-based diversity sampling performed better on average (as in Section 5.5), outperforming the BALD sampling method with dropout mechanisms and uncertainty sampling. As analysed in the rationale: in scenarios with extreme data imbalance, the reason behind its standout performance is that this strategy selects samples that cover diverse and representative information across multiple categories, providing a more comprehensive perspective for the model.

The fourth objective aimed to analyse the explanatory variable annotations. As evaluated in Section 5.6, after applying diversity sampling, using more than one high-quality explanation per iteration can significantly enhance the model’s performance while reducing response time.

In conclusion, the entire project accomplished the construction of the whole lifecycle, including design (Chapter 3), implementation (Chapter 4), and evaluation (Chapter 5). The model’s performance can be more effectively improved in the pool-based active learning framework by combining semantic-based diversity sampling and high-quality explanation annotations. The annotation process combined with the OpenAI model can serve as an alternative when manual annotation costs are high.

6.2 Limitations

First and foremost, it is worth noting that the Human-in-the-loop framework inherently has limitations in enhancing model performance. Specifically, during the annotation process, if the primary focus of users is to reduce human labour costs, the system might lean more towards using explanations generated by the OpenAI model. However, the prompt word limit of 4097 in OpenAI limits the number of samples that can undergo semantic analysis in each round [35]. As a result, explanations generated by OpenAI may not be comprehensive enough to consider all relevant information, possibly constraining the model’s performance. When considering manual annotation, there is a risk of mislabeling. Although we discussed the model’s robustness to noise in Section 5.6, excessive mislabeling could still significantly impact the model’s performance.

Further, the iterative stopping parameter of the project was determined based on the CrisisNLP dataset. For different tasks or datasets, these parameters might require adjustments. As this is not an automated adjustment process, this aspect also limits the project’s application flexibility. Moreover, while the semantic diversity sampling strategy excels in situations with severe data imbalance, its advantages might diminish if the balance of the data changes. Lastly, the time and computational resources needed for model updates must be considered. In an environment with limited resources, the high computational costs could limit the feasibility of the project.

6.3 Future Work

As this project progresses, there are multiple potential areas for further research. Addressing these areas, this section proposes the following valuable research directions:

Addressing Data Imbalance: In this study, Section 5.5 found that the active learning strategy based on semantic diversity performs superiorly in severely imbalanced data scenarios. However, given that data distribution in real-world scenarios might frequently change, future research will explore combining other strategies or integrating with uncertainty-based strategies, to adapt to changing data distributions and enhance the model’s robustness.

Optimizing Explanation Generation: Currently, the generation of explanations takes into account the capabilities of the OpenAI model and manual annotations. Nevertheless, its efficacy relies on repeated experimental validation. To systematically assess and optimize explanations, we plan to employ adversarial testing and credibility evaluations [65]. By proposing a series of target questions or misleading

queries, we aim to unveil potential weaknesses in the explanation generation process and evaluate based on the outcomes.

Broadening Framework Application: While this project primarily focuses on applying active learning to ExpBERT-based text classification tasks, considering the strategy’s efficacy, future research might think about applying this framework to a broader range of natural language processing (NLP) tasks, such as Named Entity Recognition (NER) [48] and Question Answering (QA) systems [8], further exploring its value in diverse scenarios.

Performance Testing in Real-world Scenarios: Although the evaluation methods in Chapter 5 are based on a specific experimental environment, to gain a deeper understanding of the model’s applicability in real-world contexts, there’s an anticipation of designing online systems and user interfaces for intuitive model performance evaluations. This would not only aid in comprehensively understanding the model’s performance in real scenarios but also provide prospects for its application in big data environments.

Bibliography

- [1] Charu C Aggarwal, Xiangnan Kong, Quanquan Gu, Jiawei Han, and S Yu Philip. Active learning: A survey. In *Data classification*, pages 599–634. Chapman and Hall/CRC, 2014.
- [2] Saleema Amershi, Maya Cakmak, William Bradley Knox, and Todd Kulesza. Power to the people: The role of humans in interactive machine learning. *Ai Magazine*, 35(4):105–120, 2014.
- [3] Dana Angluin. Queries and concept learning. *Machine learning*, 2:319–342, 1988.
- [4] Amos Azaria. Chatgpt usage and limitations. 2022.
- [5] Yoram Baram, Ran El Yaniv, and Kobi Luz. Online choice of active learning algorithms. *Journal of Machine Learning Research*, 5(Mar):255–291, 2004.
- [6] Christopher M Bishop and Nasser M Nasrabadi. *Pattern recognition and machine learning*, volume 4. Springer, 2006.
- [7] Michael Bloodgood. Support vector machine active learning algorithms with query-by-committee versus closest-to-hyperplane selection. In *2018 IEEE 12th International Conference on Semantic Computing (ICSC)*, pages 148–155. IEEE, 2018.
- [8] Yasaman Boreshban, Seyed Morteza Mirbostani, Gholamreza Ghassem-Sani, Seyed Abolghasem Mirroshandel, and Shahin Amiriparian. Improving question answering performance using knowledge distillation and active learning. *Engineering Applications of Artificial Intelligence*, 123:106137, 2023.
- [9] Robert Burbidge, Jem J Rowland, and Ross D King. Active learning for regression based on query by committee. In *Intelligent Data Engineering and Automated Learning-IDEAL 2007: 8th International Conference, Birmingham, UK, December 16-19, 2007. Proceedings 8*, pages 209–218. Springer, 2007.
- [10] Gavin C Cawley. Baseline methods for active learning. In *Active Learning and Experimental Design workshop In conjunction with AISTATS 2010*, pages 47–57. JMLR Workshop and Conference Proceedings, 2011.
- [11] Ido Dagan and Sean P Engelson. Committee-based sampling for training probabilistic classifiers. In *Machine Learning Proceedings 1995*, pages 150–157. Elsevier, 1995.
- [12] Tivadar Danka and Peter Horvath. modal: A modular active learning framework for python. *arXiv preprint arXiv:1805.00979*, 2018.
- [13] Pieter-Tjerk De Boer, Dirk P Kroese, Shie Mannor, and Reuven Y Rubinstein. A tutorial on the cross-entropy method. *Annals of operations research*, 134:19–67, 2005.
- [14] Rati Devidze, Farnam Mansouri, Luis Haug, Yuxin Chen, and Adish Singla. Understanding the power and limitations of teaching with imperfect knowledge. *arXiv preprint arXiv:2003.09712*, 2020.
- [15] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*, 2018.
- [16] Bo Du, Zengmao Wang, Lefei Zhang, Liangpei Zhang, Wei Liu, Jialie Shen, and Dacheng Tao. Exploring representativeness and informativeness for active learning. *IEEE transactions on cybernetics*, 47(1):14–26, 2015.
- [17] John J Dudley and Per Ola Kristensson. A review of user interface design for interactive machine learning. *ACM Transactions on Interactive Intelligent Systems (TiiS)*, 8(2):1–37, 2018.

-
- [18] Alan H Fielding and John F Bell. A review of methods for the assessment of prediction errors in conservation presence/absence models. *Environmental conservation*, 24(1):38–49, 1997.
 - [19] Yifan Fu, Xingquan Zhu, and Bin Li. A survey on instance selection for active learning. *Knowledge and information systems*, 35:249–283, 2013.
 - [20] Yarin Gal and Zoubin Ghahramani. Dropout as a bayesian approximation: Representing model uncertainty in deep learning. In *international conference on machine learning*, pages 1050–1059. PMLR, 2016.
 - [21] Yarin Gal, Riashat Islam, and Zoubin Ghahramani. Deep bayesian active learning with image data. In *International conference on machine learning*, pages 1183–1192. PMLR, 2017.
 - [22] Ethan Goan and Clinton Fookes. Bayesian neural networks: An introduction and survey. *Case Studies in Applied Bayesian Data Science: CIRM Jean-Morlet Chair, Fall 2018*, pages 45–87, 2020.
 - [23] Shantanu Godbole, Abhay Harpale, Sunita Sarawagi, and Soumen Chakrabarti. Document classification through interactive supervision of document and term labels. In *Knowledge Discovery in Databases: PKDD 2004: 8th European Conference on Principles and Practice of Knowledge Discovery in Databases, Pisa, Italy, September 20-24, 2004. Proceedings 8*, pages 185–196. Springer, 2004.
 - [24] Ian Goodfellow, Yoshua Bengio, and Aaron Courville. *Deep learning*. MIT press, 2016.
 - [25] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial networks. *Communications of the ACM*, 63(11):139–144, 2020.
 - [26] Margherita Grandini, Enrico Bagli, and Giorgio Visani. Metrics for multi-class classification: an overview. *arXiv preprint arXiv:2008.05756*, 2020.
 - [27] W Keith Hastings. Monte carlo sampling methods using markov chains and their applications. 1970.
 - [28] Geoffrey E Hinton, Nitish Srivastava, Alex Krizhevsky, Ilya Sutskever, and Ruslan R Salakhutdinov. Improving neural networks by preventing co-adaptation of feature detectors. *arXiv preprint arXiv:1207.0580*, 2012.
 - [29] Neil Houlsby, Ferenc Huszár, Zoubin Ghahramani, and Máté Lengyel. Bayesian active learning for classification and preference learning. *arXiv preprint arXiv:1112.5745*, 2011.
 - [30] Rong Hu, Brian Mac Namee, and Sarah Jane Delany. Active learning for text classification with reusability. *Expert systems with applications*, 45:438–449, 2016.
 - [31] Sagar Imambi, Kolla Bhanu Prakash, and GR Kanagachidambaresan. Pytorch. *Programming with TensorFlow: Solution for Edge Computing Applications*, pages 87–104, 2021.
 - [32] Muhammad Imran, Prasenjit Mitra, and Carlos Castillo. Twitter as a lifeline: Human-annotated twitter corpora for nlp of crisis-related messages. *arXiv preprint arXiv:1605.05894*, 2016.
 - [33] Akhil Kadiyala and Ashok Kumar. Applications of python to evaluate environmental data science problems. *Environmental Progress & Sustainable Energy*, 36(6):1580–1586, 2017.
 - [34] Ross D King, Kenneth E Whelan, Ffion M Jones, Philip GK Reiser, Christopher H Bryant, Stephen H Muggleton, Douglas B Kell, and Stephen G Oliver. Functional genomic hypothesis generation and experimentation by a robot scientist. *Nature*, 427(6971):247–252, 2004.
 - [35] Ketu Korini and Christian Bizer. Column type annotation using chatgpt. *arXiv preprint arXiv:2306.00745*, 2023.
 - [36] Agustinus Kristiadi, Matthias Hein, and Philipp Hennig. Being bayesian, even just a bit, fixes overconfidence in relu networks. In *International conference on machine learning*, pages 5436–5446. PMLR, 2020.
 - [37] K Lang and E Baum. Query learning can work poorly when a human oracle is used. *ieee intl. In Joint Conference on Neural Networks*, 1992.
-

- [38] Jeannette Lawrence. *Introduction to neural networks*. California Scientific Software, 1993.
- [39] Florian Laws and Hinrich Schütze. Stopping criteria for active learning of named entity recognition. In *Proceedings of the 22nd International Conference on Computational Linguistics (Coling 2008)*, pages 465–472, 2008.
- [40] David D Lewis. A sequential algorithm for training text classifiers: Corrigendum and additional data. In *Acm Sigir Forum*, volume 29, pages 13–19. ACM New York, NY, USA, 1995.
- [41] Tianyang Lin, Yuxin Wang, Xiangyang Liu, and Xipeng Qiu. A survey of transformers. *AI Open*, 2022.
- [42] Bill MacCartney. *Natural language inference*. Stanford University, 2009.
- [43] David JC MacKay. A practical bayesian framework for backpropagation networks. *Neural computation*, 4(3):448–472, 1992.
- [44] Andrew McCallum, Kamal Nigam, et al. Employing em and pool-based active learning for text classification. In *ICML*, volume 98, pages 350–358. Citeseer, 1998.
- [45] Richard McCreadie, Cody Buntain, and Ian Soboroff. TREC incident streams: Finding actionable information on social media. 2019.
- [46] Wes McKinney et al. pandas: a foundational python library for data analysis and statistics. *Python for high performance and scientific computing*, 14(9):1–9, 2011.
- [47] Chris J Michael, Dina Acklin, and Jaelle Scheuerman. On interactive machine learning and the potential of cognitive feedback. *arXiv preprint arXiv:2003.10365*, 2020.
- [48] Behrang Mohit. Named entity recognition. In *Natural language processing of semitic languages*, pages 221–245. Springer, 2014.
- [49] Robert Munro Monarch. *Human-in-the-Loop Machine Learning: Active learning and annotation for human-centered AI*. Simon and Schuster, 2021.
- [50] Eduardo Mosqueira-Rey, Elena Hernández-Pereira, David Alonso-Ríos, José Bobes-Bascarán, and Ángel Fernández-Leal. Human-in-the-loop machine learning: A state of the art. *Artificial Intelligence Review*, 56(4):3005–3054, 2023.
- [51] Vikram Mullachery, Aniruddh Khera, and Amir Husain. Bayesian neural networks. *arXiv preprint arXiv:1801.07710*, 2018.
- [52] Shikhar Murty, Pang Wei Koh, and Percy Liang. Expbert: Representation engineering with natural language explanations. *arXiv preprint arXiv:2005.01932*, 2020.
- [53] Tomoyuki Myojin, Shintaro Hashimoto, and Naoki Ishihama. Detecting uncertain bnn outputs on fpga using monte carlo dropout sampling. In *Artificial Neural Networks and Machine Learning–ICANN 2020: 29th International Conference on Artificial Neural Networks, Bratislava, Slovakia, September 15–18, 2020, Proceedings, Part II 29*, pages 27–38. Springer, 2020.
- [54] Radford M Neal. *Bayesian learning for neural networks*, volume 118. Springer Science & Business Media, 2012.
- [55] Travis E Oliphant et al. *Guide to numpy*, volume 1. Trelgol Publishing USA, 2006.
- [56] Bo Pang, Erik Nijkamp, and Ying Nian Wu. Deep learning with tensorflow: A review. *Journal of Educational and Behavioral Statistics*, 45(2):227–248, 2020.
- [57] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, et al. Pytorch: An imperative style, high-performance deep learning library. *Advances in neural information processing systems*, 32, 2019.
- [58] Hao Peng, Shixin Guo, Dandan Zhao, Yiming Wu, Jianming Han, Zhe Wang, Shouling Ji, and Ming Zhong. Query-efficient model extraction for text classification model in a hard label setting. *Journal of King Saud University-Computer and Information Sciences*, 35(4):10–20, 2023.

-
- [59] Pietro Portolani, Alessandro Brusafferri, Andrea Ballarino, and Matteo Matteucci. Uncertainty in predictive process monitoring. In *International Conference on Information Processing and Management of Uncertainty in Knowledge-Based Systems*, pages 547–559. Springer, 2022.
 - [60] Sumanth Prabhu, Moosa Mohamed, and Hemant Misra. Multi-class text classification using bert-based active learning. *arXiv preprint arXiv:2104.14289*, 2021.
 - [61] Alessandro Prest, Cordelia Schmid, and Vittorio Ferrari. Weakly supervised learning of interactions between humans and objects. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 34(3):601–614, 2011.
 - [62] Hou-Duo Qi and Xiaoming Yuan. Computing the nearest euclidean distance matrix with low embedding dimensions. *Mathematical programming*, 147(1-2):351–389, 2014.
 - [63] Gonzalo Ramos, Christopher Meek, Patrice Simard, Jina Suh, and Soroush Ghorashi. Interactive machine teaching: a human-centered approach to building machine-learned models. *Human-Computer Interaction*, 35(5-6):413–451, 2020.
 - [64] Pengzhen Ren, Yun Xiao, Xiaojun Chang, Po-Yao Huang, Zhihui Li, Brij B Gupta, Xiaojiang Chen, and Xin Wang. A survey of deep active learning. *ACM computing surveys (CSUR)*, 54(9):1–40, 2021.
 - [65] 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*, pages 1135–1144, 2016.
 - [66] Christopher Schröder, Lydia Müller, Andreas Niekler, and Martin Potthast. Small-text: Active learning for text classification in python. *arXiv preprint arXiv:2107.10314*, 2021.
 - [67] Christopher Schröder and Andreas Niekler. A survey of active learning for text classification using deep neural networks. *arXiv preprint arXiv:2008.07267*, 2020.
 - [68] Ozan Sener and Silvio Savarese. Active learning for convolutional neural networks: A core-set approach. *arXiv preprint arXiv:1708.00489*, 2017.
 - [69] Burr Settles. Active learning literature survey. 2009.
 - [70] H Sebastian Seung, Manfred Opper, and Haim Sompolinsky. Query by committee. In *Proceedings of the fifth annual workshop on Computational learning theory*, pages 287–294, 1992.
 - [71] Sagar Sharma, Simone Sharma, and Anidhya Athaiya. Activation functions in neural networks. *Towards Data Sci*, 6(12):310–316, 2017.
 - [72] Yucheng Shi, Hehuan Ma, Wenliang Zhong, Gengchen Mai, Xiang Li, Tianming Liu, and Junzhou Huang. Chatgraph: Interpretable text classification by converting chatgpt knowledge to graphs. *arXiv preprint arXiv:2305.03513*, 2023.
 - [73] Nitish Srivastava, Geoffrey Hinton, Alex Krizhevsky, Ilya Sutskever, and Ruslan Salakhutdinov. Dropout: a simple way to prevent neural networks from overfitting. *The journal of machine learning research*, 15(1):1929–1958, 2014.
 - [74] John E Tomaszewski. Overview of the role of artificial intelligence in pathology: the computer as a pathology digital assistant. In *Artificial intelligence and deep learning in pathology*, pages 237–262. Elsevier, 2021.
 - [75] Evgenii Tsymbalov, Maxim Panov, and Alexander Shapeev. Dropout-based active learning for regression. In *Analysis of Images, Social Networks and Texts: 7th International Conference, AIST 2018, Moscow, Russia, July 5–7, 2018, Revised Selected Papers 7*, pages 247–258. Springer, 2018.
 - [76] Devis Tuia, Michele Volpi, Loris Copa, Mikhail Kanevski, and Jordi Munoz-Mari. A survey of active learning algorithms for supervised remote sensing image classification. *IEEE Journal of Selected Topics in Signal Processing*, 5(3):606–617, 2011.
-

- [77] Adina Williams, Nikita Nangia, and Samuel Bowman. A broad-coverage challenge corpus for sentence understanding through inference. In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)*, pages 1112–1122, New Orleans, Louisiana, June 2018. Association for Computational Linguistics.
- [78] Xufeng Zhao, Mengdi Li, Cornelius Weber, Muhammad Burhan Hafez, and Stefan Wermter. Chat with the environment: Interactive multimodal perception using large language models. *arXiv preprint arXiv:2303.08268*, 2023.

Appendix A

Github Repository

Github Repository: [https://github.com/sxy1813082/Data-Science-Project/ExpBERT lab and thesis/BERT-IE-main/BERT-IE-main](https://github.com/sxy1813082/Data-Science-Project/ExpBERT%20lab%20and%20thesis/BERT-IE-main/BERT-IE-main)