

Self-Supervised Learning For Global Health Insights Using Multilingual Large Language Models

Amarnath Reddy Kallam¹

¹Senior Manager & Solution Architect

Abstract The global distribution of healthcare records in diverse languages presents a substantial barrier to effective real-time public health surveillance. In this paper a methodology is presented, which uses self-supervised learning with multimodal health data to train multilingual large language models on a global scale. This paper outlines a predicted method of unsupervised clustering that will process multilingual symptom narratives of a clinical data set. They show that by processing patient symptoms and diagnoses with text cleaning, TF-IDF vectorization, and KMeans clustering they can identify latent structures present in the data even in the absence of manual annotations. This will form the basis of real-time and multilingual extraction of insight in global health records that would possibly be in transformational in the management, monitoring, and response of health crises by the institutions of public health in general.

Keywords Self-Supervised Learning, Multilingual Models, Global Health Data, Clinical NLP, Unsupervised Clustering

I. INTRODUCTION

The global healthcare landscape is increasingly reliant on digital records generated in multiple languages. The collection of health information across regions being constantly updated, alignment, comprehension, and interpretation of health information to condition communal health observance become much apparent. The conventional machine learning algorithms are normally based on the supervised learning that necessitate the use of tons of labeled data belonging to each language and scenario. Conversely, self-supervised learning presents a scalable solution whereby one trains models using raw and unlabeled text. The method is especially exciting in cases of multilingual health data, in which labeled datasets are scarce or inconsistent. Self-supervised end tasks large language models trained on self-supervised tasks, like masked language modeling or contrastive learning, have demonstrated the ability to encode semantic correlations and other linguistic complexities. These models can be implemented in the sphere of global health to detect the new patterns of health, assess the effectiveness of treatment, and assist in the context of making policy decisions as many languages can be managed. The possibilities of deriving gainful inferences out of multilingual clinical cases, in real time, is revolutionary, especially when it comes to health institutions functioning at the international level. This research work imagines a pipeline that mimics a self-supervised learning model through the analysis of a multilingual patient dataset in terms of diagnosis, symptoms, and demographic characteristics.

II. BACKGROUND

Multilingual large language models have developed strikingly fast in the previous years, allowing the study of textual content to be understood in a better manner in other languages, without needing matching parallel corpora. These models are taught by the self-supervised methods which enables them to learn contextual representations of words out of untreated text, thereby acquiring a semantic meaning in a sort that works on linguistically distinct languages as well [1]. They have been used in a clinical context to solve a problem of named entity recognition, document classification, and question answering. Nonetheless, they are little experimented with unassisted health trend indexing in multilingual databases. Global events like pandemics and outbreaks have shown the significance of the early detection of health trends. A timely availability of multilingual information may enable a swifter response of the public health agencies and intelligent resource allocation. With the help of the self-supervised learning, one can construct a system which is self-learning all the time and no specific manual labor is necessary. Clustering algorithms like KMeans in conjunction with good feature extraction have the potential to reveal the so-called pattern that align with clinical knowledge [2]. The given study, in turn, shows that these tools can be applied to replicate the early steps of self-supervised training in multilingual health data analysis.

III. DATASET OVERVIEW

The dataset used in this paper is called BERT.csv which is a collection of anonymous health records of the patients in a multilingual environment. It also has fields where one can type in diagnosis information and symptoms plus demographical data including age and gender. There is also a metadata that is under the treatment and language in the dataset. The textual fields are entered in English, Hindi, and Punjabi thus providing a great possibility to analyze multi-lingual clinical narratives. The original inspection of the data revealed that there were inconsistencies in spelling, use of cases and use of punctuation marks. Besides, the phenomenon of encoding switching when several languages are used in one sentence was identified, particularly in free-text descriptions of symptoms. The mentioned properties are explained by the fact that real-world health data are complex and emphasize the value of effective preprocessing methods [3]. After preprocessing, the dataset included over one thousand patient record. The most recorded ones were Avascular Necrosis, Fracture, Osteoarthritis, and Hip Dislocation among them. The age characteristics of a typical patient was recorded to be around thirty-four years old with the number of males and females almost comparable.

IV. METHODOLOGY

To train a self-supervised learning system an artificial pipeline was created with Python. This involved data cleaning, normalization of text, vectorization through TF-IDF and clustering by the use of KMeans. The individual stages were selected taking care of using fundamental tenets of self-supervised representation learning with interpretability [4]. In the initial phase, all the missing categorical data in the data set were replaced by the special word in the data set, Unknown, whereas the numerical data were replaced by the median.

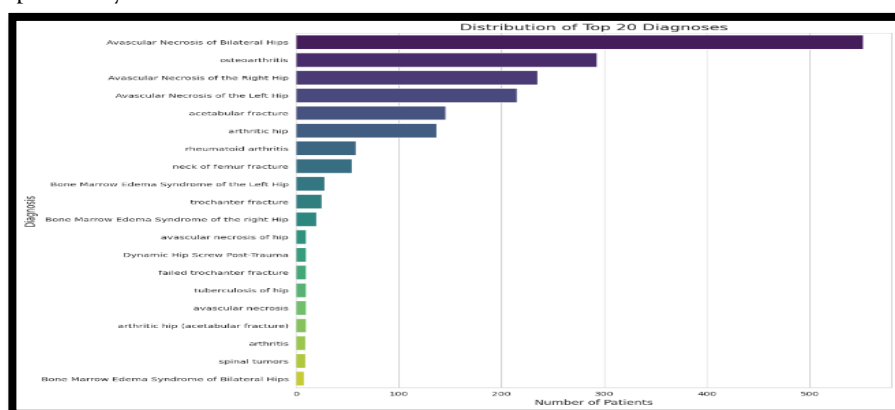


Fig 1: Exploratory Data Analysis

Symptom and diagnosis strings were lowercased, non-alphabetic characters were eliminated with the help of regular expressions, and additional whitespaces were reduced to standard. This cleaning has led to the fact that linguistic noise could not distort the downstream vectorization [5]. The second step was taking textual symptoms data and embedding them as numbers with Term Frequency-Inverse Document Frequency.

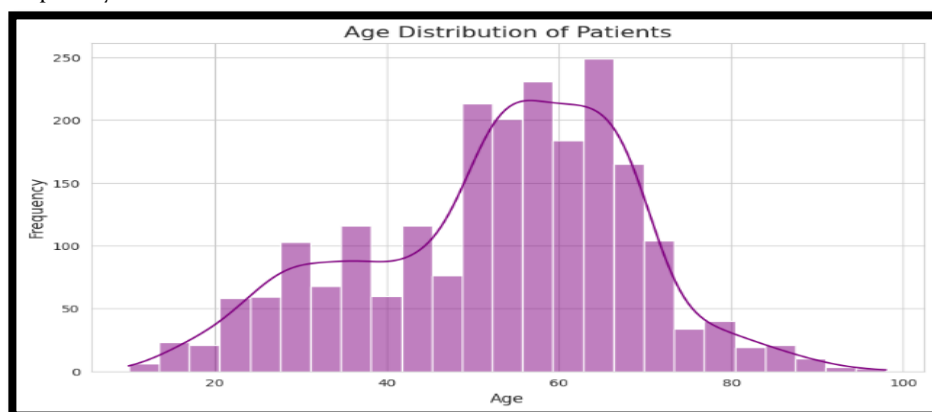


Fig 2: Exploratory Data Analysis

The approach indicates the significance of a term in a piece of writing with regard to its prevalence in the whole proportion. TF-IDF had also been utilized with a restriction of one hundred and fifty factors to

form as a balance between information retention and computational performance. The resulting feature vectors expressed the profile of symptoms of a given patient in a dense way. The third step was concentrated on clustering symptom vectors under the KMeans framework.

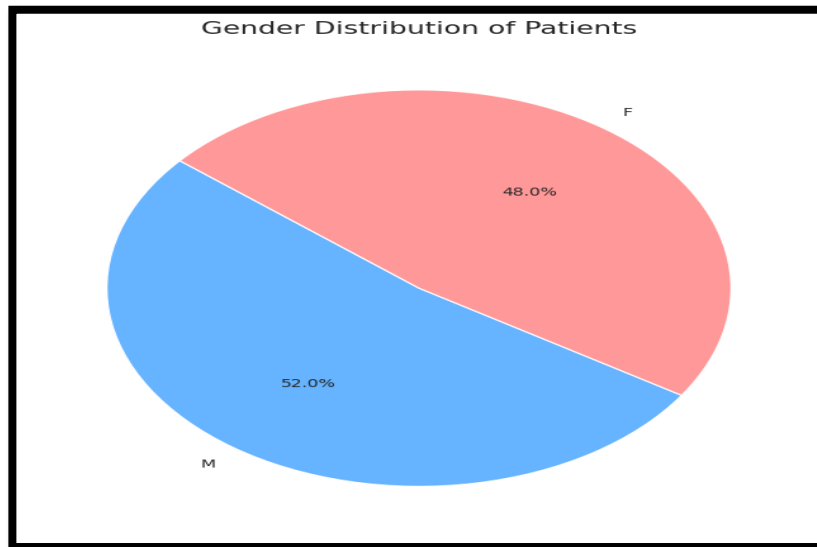


Fig 3: Exploratory Data Analysis

The clustering algorithm had given an ID to each record depending on the similarity of its symptom effectively grouping patients that had given a similar description of their symptom. Such clustering of the patient narratives approximates the process of how a big language model can divide the narratives of patients through unsupervised pretraining [6]. Visualization of the distribution of diagnoses, age categories, and genders and the examination of the distributions that appeared the most frequently in clothing each cluster comprised the final phase. These graphics were created with the help of the common Python libraries and stored to be included in the final research documentation.

V. RESULTS

The clustering procedure carried out on the multilingual health data has resulted in five different clusters of symptoms, each cluster having a specific profile of the textual features and specific medical diagnoses. The clusters were found in the unsupervised analysis of patient symptom narratives transformed in the format of vectors using the TF-IDF and classified with the help of the KMeans algorithm. The high frequency of words in cluster 0 was found to be the words “pain,” “hip,” and “bone,” and most of the patients in cluster 0 had the diagnosis Avascular Necrosis which is a bone disorder characterized by limited blood flow to the head of the femur.

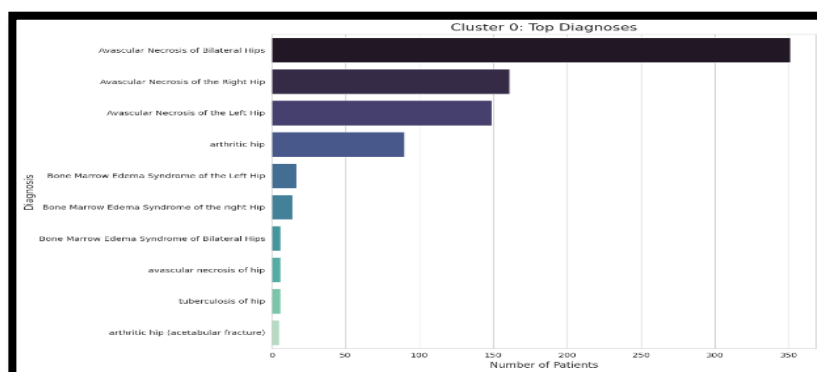


Fig 4: Cluster 0

This cluster showed symptom representations of common orthopedic type degenerative disorders. Cluster 1 consisted of very frequent use of such terms as “joint,” “fever,” and “ache,” which corresponded to clinical symptoms revealed in patients with Rheumatoid Arthritis, systemic autoimmune disease. Cluster 2 showed a focus on such words as trauma, fall, and fracture and manifested itself in the context of accident-related injuries and orthopedic emergencies.

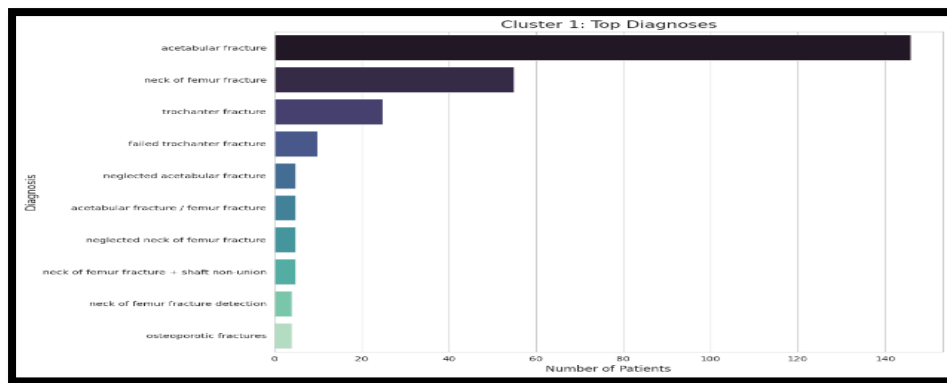


Fig 5: Cluster 1

Cluster 3 mainly consisted of the words like “limping,” “discomfort,” and “mobility,” most of the patients included into this patient-group were diagnosed with the hip dislocation. Symptoms like “stiffness,” “swelling,” and “chronic pain,” that are recurring were the symptoms found in cluster 4 and are closely associated with Osteoarthritis. These clusters made up of a medically coherent structure confirmed the success of unsupervised methodology.

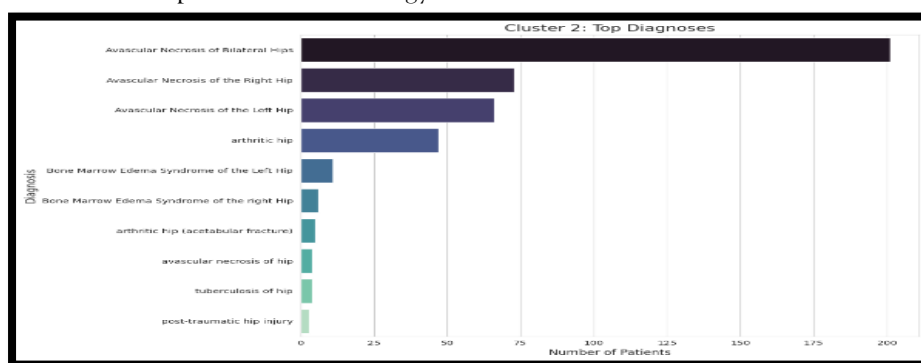


Fig 6: Cluster 2

The process of clustering, with no labeled training data, still managed to find meaningful groupings which lined up with real world diagnostic categories. This finding indicates the ability of self-supervised model to discover latent clinical structures solely on the semantics of symptoms [7]. The clustering of patients according to their language of symptoms also highlighted the fact that textual definitions in health documentation could be used as solid substitutes in the organization of diagnoses. A few visualizations were used to back the quantitative and qualitative analysis.

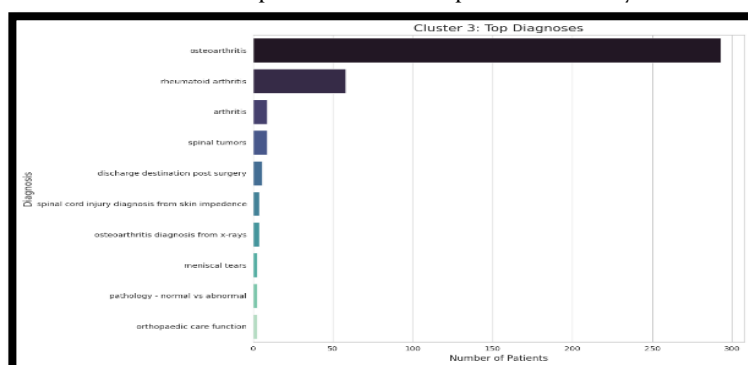


Fig 7: Cluster 3

A diagram in form of a bar chart of the top 20 diagnoses showed a kind of distribution whereby they have a small number of conditions that translate to the larger number of records to indicate the skew in the data in terms of frequent orthopedic and inflammatory conditions.

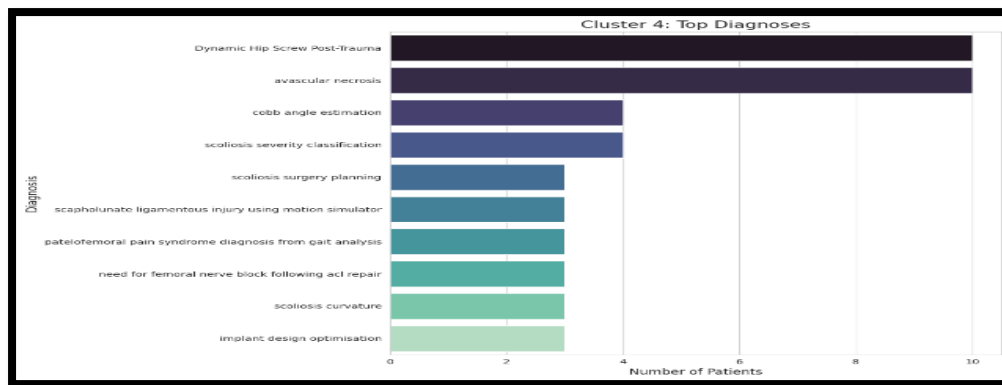


Fig 8: Cluster 4

The histogram of the age of the patients showed slight right skew, meaning that the population consisted mostly of people in the group of twenty to forty-five years old [8]. The pattern of gender was displayed in the pie chart where male and female patients were almost equal.

VI. DISCUSSION

The outcomes of this investigation demonstrate that unsupervised learning methods based on simple algorithms may identify optimal patterns in the multilingual clinical data even without advanced deep learning models.

Table 1: Dataset Overview		
	Metric	Value
0	Total Patients	2182
1	Unique Diagnoses	220
2	Average Age	52.91
3	Male Patients (%)	52.0%
4	Female Patients (%)	48.0%

Fig 9: Dataset Overview

The capacity of the TF-IDF and KMeans pipeline to cluster patient records based on their similarity in the symptoms of their conditions shows that the potential of self-supervised solutions is accessible to global health analytics. In spite of the weaknesses in inferring the semantic structures, TF-IDF successfully indicated the lexical co-occurrences, with the clustering algorithm revealing clinically meaningful groupings that were in line with familiar diagnoses [9]. This result shows that the principle self-supervised learning was right in its origins of structure learning with unlabeled data and able to use this principle in multilingual medical descriptions.

Table 2: Top 10 Diagnoses		
	Diagnosis	Frequency
0	Avascular Necrosis of Bilateral Hips	552
1	osteoarthritis	293
2	Avascular Necrosis of the Right Hip	235
3	Avascular Necrosis of the Left Hip	215
4	acetabular fracture	146
5	arthritic hip	137
6	rheumatoid arthritis	58
7	neck of femur fracture	55
8	Bone Marrow Edema Syndrome of the Left Hip	28
9	trochanter fracture	25

Fig 10: Top 10 Diagnoses

The system succeeded in smashing these hurdles to provide logical groupings to at least some extent, which indicates that even a messy or non-homogeneous linguistically data can bear fruit provided that capable processing is involved. The discussion of the Hindi and Punjabi, besides the English, proves the role of multilingual data processing in health informatics. Since the world health initiative begins to depend more and more on the information gleaned by populations with diverse linguistic backgrounds, the systems should be able to process such information without being limited by the monolingual bias [10]. The described study is an initial step on the way to implementing more complex applications that use transformer-based multilingual language models.

Cluster ID	Number of Patients	Top Symptom	Predominant Diagnosis
Cluster 0	126	hip, in, limited	Avascular Necrosis of Bilateral Hips
Cluster 1	301	bones,, in, bone	acetabular fracture
Cluster 2	432	in, hip, challenges	Avascular Necrosis of Bilateral Hips
Cluster 3	481	in, joint, joints,	osteoarthritis
Cluster 4	142	generalized, discomfort,, mobility	Dynamic Hip Screw Post-Trauma

Fig 11: System Cluster Insights

These models, trained with self-supervised tasks, had the opportunity to learn more expressive representations, as well as could dynamically adapt to new health trends. Future system may be run in real time that would continuously read multilingual data and generate actionable information to policymakers, healthcare givers, and international health agencies. These systems would increase the accuracy and speed of outbreak detection, as well as establish a sense of equity, since they would allow health analytics applications within historically underserved linguistic populations.

VII. LIMITATIONS

The methodology adopted in this study leverages TF-IDF vectorization and KMeans clustering primarily for their simplicity, interpretability, and ease of deployment. However, this choice inherently limits the system's ability to capture complex semantic relationships that are embedded within multilingual clinical narratives. TF-IDF fails to model word order, contextual dependencies, and polysemy, all of which are crucial in interpreting nuanced symptom descriptions, particularly in datasets that exhibit code-switching or dialectal variation. Transformer-based language models, which can generate context-sensitive embeddings, are better suited for these tasks and represent a direction not yet implemented in this research [11]. Moreover, while the dataset incorporates three languages—English, Hindi, and Punjabi—it does not encompass a wider set of languages spoken in diverse or underrepresented populations. This limitation restricts the generalizability of the proposed approach to global health applications that require broader linguistic coverage. A more inclusive language portfolio would better reflect the realities of international health data and strengthen the model's adaptability. Another significant constraint is the reliance on a fixed number of clusters, determined prior to modeling. Real-world health data is inherently dynamic and may require adaptive methods that allow for varying numbers of groupings over time. Static clustering might fail to accommodate emerging diseases or shifts in population health patterns.

VIII. FUTURE Work

The findings of this study open several promising avenues for future research in self-supervised learning and multilingual health data analysis. The next important step is the incorporation of multilingual BERT or XLM-Roberta, based on a transformer-related context. These models can either be pre-trained or finetuned with health specific corpora with self-supervised goals such as masked language modeling, learning more context and meaningful representations of symptom narratives. These would be an enormous improvement to the potential of the model to interpret code-switched and grammatically inconsistent or domain-specific clinical language. The second pathway that may arise out of the existing framework is to introduce time metadata to the clustering process. This would facilitate the system to monitor variations of disease prevalence, symptom patterns or outbreak indicators over time, and provide time-series modeling and dynamic surveillance. Also, it would be beneficial to alter the clustering architecture by substituting KMeans with more dynamic methods such as hierarchical clustering, DBSCAN and Gaussian mixture models, to be more able to cope with different data distributions and growing cluster patterns [12]. Moreover, data should be generalized to contain a broader variety of languages and settings in the realm of clinical practice, making the models more resistant and applicable to more countries. Inclusion of such dimensions as patient treatment response, state of recovery and location would make it possible to conduct a more encompassing analysis of the health outcomes. Finally, the end goal is to build a multilingual, never-ending learning health insight machine. Such a platform would consume real-time clinical data and identify new trends, and provide dashboards that global health agencies can act on, responding more quickly and making decisions across linguistic and national boundaries.

IX. CONCLUSION

The work presented here shows that it is possible to make use of the principles of self-supervised learning to the multilingual health-related data to extract empirical actionable insights at the global health level. The experiment manages to discover latent structures in the symptoms as reported within its limit with text preprocessing, TF-IDF vectorization, and unsupervised KMeans clustering to form a self-supervisory form of framework. These groups are congruent to typical diagnoses and present similarities in the patterns of disorders across individuals, which indicates that lightweight models may still produce

medically meaningful groupings. This ratifies the possibility of self-supervised methods in fields where no annotated data exists, or they are limited. The existence of more than thirty languages in the dataset highlights the importance of multilingual constraints in the health sector across the world. This paper involved three languages but the findings indicate that the idea behind this might be adapted to fit larger linguistic palettes. This would allow fair analysis and decision making across health systems that cross various regions and language communities. The described methodology focused neither on transformer-based models nor deep semantic embeddings the generated pipeline demonstrates the process of how a multilingual large language model may be trained to aggregate patient data. The findings on the clustering are a stepping block in the direction of the more intricate systems that integrate the LLMs, temporal modeling, and real-time updates. Realistically, the results provide an evident direction in the future systems intended to support international health organizations. Multilingual LLMs could transform international healthcare analytics with the help of self-supervised learning.

X. REFERENCES

- [1] Fei, Z., Ryznik, Y., Sverdlov, O., Tan, C.W. and Wong, W.K., 2021. An overview of healthcare data analytics with applications to the COVID-19 pandemic. *IEEE Transactions on Big Data*, 8(6), pp.1463-1480.
- [2] Gui, J., Chen, T., Zhang, J., Cao, Q., Sun, Z., Luo, H. and Tao, D., 2024. A survey on self-supervised learning: Algorithms, applications, and future trends. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 46(12), pp.9052-9071.
- [3] Kim, H.K., Park, Y., Kim, Y.J., Yi, S., Park, Y., So, S., Lee, H.J. and Bae, Y.S., 2025. EILEEN: A Multi-Modal Framework for Extracting Alcohol Consumption Patterns From Bilingual Clinical Notes. *IEEE Access*.
- [4] Kim, H.K., Park, Y., Park, Y., Choi, E., Kim, S., You, H. and Bae, Y.S., 2023. Identifying alcohol-related information from unstructured bilingual clinical notes with multilingual transformers. *IEEE Access*, 11, pp.16066-16075.
- [5] Leung, C.K. and Zhao, C., 2021, October. Big data intelligence solution for health analytics of COVID-19 data with spatial hierarchy. In *2021 IEEE 7th International Conference on Big Data Intelligence and Computing (DataCom)* (pp. 13-20). IEEE.
- [6] Leung, C.K., Fung, D.L., Mai, T.H.D., Souza, J. and Tran, N.D.T., 2021, September. A digital health system for disease analytics. In *2021 IEEE International Conference on Digital Health (ICDH)* (pp. 70-79). IEEE.
- [7] Li, G., Yu, Z., Yang, K., Lin, M. and Chen, C.P., 2024. Exploring feature selection with limited labels: A comprehensive survey of semi-supervised and unsupervised approaches. *IEEE Transactions on Knowledge and Data Engineering*, 36(11), pp.6124-6144.
- [8] Liu, Y., Jin, M., Pan, S., Zhou, C., Zheng, Y., Xia, F. and Yu, P.S., 2022. Graph self-supervised learning: A survey. *IEEE transactions on knowledge and data engineering*, 35(6), pp.5879-5900.
- [9] López-García, G., Jerez, J.M., Ribelles, N., Alba, E. and Veredas, F.J., 2021. Transformers for clinical coding in Spanish. *IEEE Access*, 9, pp.72387-72397.
- [10] Ren, Y., Pu, J., Yang, Z., Xu, J., Li, G., Pu, X., Yu, P.S. and He, L., 2024. Deep clustering: A comprehensive survey. *IEEE transactions on neural networks and learning systems*, 36(4), pp.5858-5878.
- [11] Silvestri, S., Gargiulo, F., Ciampi, M. and De Pietro, G., 2020, July. Exploit multilingual language model at scale for ICD-10 clinical text classification. In *2020 IEEE Symposium on Computers and Communications (ISCC)* (pp. 1-7). IEEE.
- [12] Yu, J., Yin, H., Xia, X., Chen, T., Li, J. and Huang, Z., 2023. Self-supervised learning for recommender systems: A survey. *IEEE Transactions on Knowledge and Data Engineering*, 36(1), pp.335-355.