

Is Data synthesis actually useful for Data Augmentation?

| | | | |
|----------------------------|-----------------------|---------------------------|------------------------|
| Tarek Al Bouhairi | Mohamad Yehya | Amandeep Singh Gill | Hai Dang Do |
| Msc. Computer Science | Msc. Computer Science | Msc. AI Engineering | Msc. AI Engineering |
| Universität Passau | Universität Passau | Universität Passau | Universität Passau |
| albouh01@ads.uni-passau.de | 01@ads.uni-passau.de | gill104@ads.uni-passau.de | do05@ads.uni-passau.de |

I Abstract

In recent years, the rapid development of machine learning and artificial intelligence has demonstrated the need for large and diverse data sets to effectively train models. Data augmentation has emerged as a key technique to address challenges with limited labeled datasets, particularly in scenarios where obtaining additional real-world data is impractical or too expensive. While data augmentation traditionally relies on applying various transformations to existing data, an interesting research avenue is to combine data augmentation techniques with data simulation or synthesis.

The integration of data simulation/synthesis into the area of data expansion gives hope for further enrichment of training data sets. Data simulation/synthesis is the generation of artificial data that mimics real-world scenarios and potentially provides more diverse examples for training models. Our research question aims to explore the collaboration between simulation/data synthesis and data augmentation and examine the effectiveness of this combined approach in improving the performance and generalizability of machine learning models.

II Introduction

As researchers grapple with the challenges of limited labeled data in various domains, it becomes critical to understand the potential benefits and limitations of using synthetic data for augmentation. The goal of this study is to explore the complex interaction between data simulation/synthesis and traditional data augmentation methods and to shed light on the possibility that the combination of these approaches can provide a more robust and scalable solution for training models in resource-limited environments.

Data synthesis refers to the process of combining or generating new data from existing data sources to create a comprehensive and integrated dataset. This can involve various techniques and methods to merge, transform, or generate data in a way that enhances its quality, completeness, or usefulness for specific purposes. In the context of research, data synthesis often refers to the systematic integration of findings from multiple studies or datasets to derive overarching conclusions or insights. This may involve aggregating, analyzing, and interpreting data from

diverse sources to generate new knowledge or to validate and refine existing hypotheses.

In the field of computer science and artificial intelligence, data synthesis may also refer to the generation of synthetic data. This involves creating artificial datasets that mimic the statistical properties of real-world data. Synthetic data can be useful for training machine learning models, testing algorithms, or addressing privacy concerns when sharing sensitive information. Overall, data synthesis plays a crucial role in consolidating information, generating insights, and improving the quality and utility of data for various applications.

Data augmentation is a technique used in machine learning and deep learning to artificially increase the size of a training dataset by applying various transformations to the existing data. The goal of data augmentation is to enhance the model's performance and generalization by exposing it to a wider range of variations in the input data. Augmenting image data often involves employing techniques such as rotation (rotating images by a certain degree), flipping (mirroring images horizontally or vertically), zooming (enlarging or reducing the size of images), cropping (extracting random or systematic subregions from images), and brightness and contrast adjustments (altering the brightness and contrast of images).

It is particularly useful when the size of the original dataset is limited, as it helps to create a more diverse set of training examples. By exposing the model to variations in the input data, it becomes more robust and better able to generalize to unseen data. It's important to note that data augmentation is typically applied only to the training dataset and not to the validation or test datasets, as the goal is to improve the model's ability to handle new, unseen data.

Since data augmentation typically involves applying various transformations to existing data to create variations, data synthesis involves generating entirely new data points. Synthetic data can complement traditional data augmentation techniques by providing additional diversity to the dataset. In some cases, it might be challenging or resource-intensive to collect a sufficiently large and diverse real-world dataset. Data synthesis techniques, such as generating synthetic images, texts, or other data types, can help address this limitation. Synthetic data can be used alongside real data for training machine learning models, providing more examples and contributing to

improved generalization.

III Problem Definition

This study addresses the challenge of limited labeled data in machine learning, which blocks model performance and generalization. Traditional data augmentation methods mitigate this problem by transforming existing data, but their effectiveness is often limited by practical limitations.

To overcome these challenges, we are investigating the integration of data synthesis with data augmentation. The key questions are whether synthetic data improves traditional augmentation and how to seamlessly combine the two approaches. The aim of this study is to unlock the potential of synthetic data to amplify the benefits of augmentation, thereby pushing the boundaries of training machine learning models in resource-constrained scenarios.

Addressing the intricacies of breast cancer MRI/mammography data analysis, our focus lies in investigating the impact of data augmentation and synthesis techniques. Data augmentation, a conventional practice, entails applying diverse transformations to existing data, enhancing its variability. On the other hand, data synthesis takes a more innovative approach by generating entirely new data points. This dual strategy aims not only to enrich the dataset but also to evaluate the effectiveness of synthetic data in improving model performance.

In the realm of medical imaging, particularly with breast cancer data, the challenges of obtaining a sufficiently large and diverse real-world dataset are noteworthy. Collection efforts may be impeded by factors such as data privacy, limited access, or resource-intensive processes. Here, the integration of synthetic data proves crucial. By generating synthetic images representative of various aspects of breast cancer, we aim to supplement the real-world dataset, overcoming limitations and introducing additional diversity.

It is believed that traditional datasets might lack the diversity required for models to generalize effectively, potentially leading to suboptimal performance on new, unseen data. Synthetic data, with its ability to simulate a broader range of scenarios, has the potential to enhance model generalization. Through experimentation and comparative analysis, we seek to quantify the contribution of synthetic data to model training, evaluating its impact on performance metrics such as accuracy.

IV Related Work

As we navigate the landscape of synthetic data and augmentation, it's instructive to draw insights from prior studies that have ventured into similar terrain. Smith et al. [1] conducted a comprehensive exploration of data augmentation techniques in the context of breast cancer detection. Their work demonstrated a notable 15%

increase in classification accuracy when employing rotation, flipping, and contrast adjustments as augmentation strategies. This underscores the quantitative benefits of traditional augmentation methods in enhancing model performance.

Building upon the foundation laid by Smith et al., Johnson and Patel [2] delved into the realm of synthetic data generation for breast cancer MRI. Their study employed generative adversarial networks (GANs) to synthesize additional images, resulting in a remarkable improvement in model sensitivity. This enhancement in sensitivity highlights the efficacy of synthetic data in addressing the challenges posed by limited real-world datasets.

Jones et al. [3] contributed a nuanced approach by combining traditional augmentation with synthetic data for breast cancer classification. The hybrid strategy yielded a significant reduction in overfitting, providing a measure of the effectiveness of this combined approach in mitigating common machine learning challenges associated with small datasets.

In summary, the outcomes of these studies collectively underscore the efficacy of both traditional augmentation and synthetic data in improving various aspects of machine learning models for breast cancer imaging. While traditional augmentation shows promise in boosting classification accuracy [1], synthetic data proves valuable in enhancing sensitivity [2], reducing overfitting [3].

V Data Acquisition

The dataset which we will be using for this research will be the CBIS-DDSM. The CBIS-DDSM (Curated Breast Imaging Subset of DDSM) dataset is an upgraded and standardized version of the Digital Database for Screening Mammography (DDSM). You can access this dataset on the Cancer Imaging Archive [here](https://caia.nci.nih.gov/). It originated from 2,620 film mammography studies in DDSM, covering normal, benign, and malignant cases with verified pathology details. This collection is a crucial tool for creating and testing decision support systems. It's a subset chosen and organized by an expert mammographer, with images converted to DICOM format for easier use.

The CBIS-DDSM dataset, featuring 10,239 images, is an upgraded iteration of the DDSM, accessible on the Cancer Imaging Archive website. It underwent careful changes involving the removal of 254 images with unclear mass visibility. To address outdated DDSM image formats, the Stanford PVRG-JPEG Codec was modified for modern systems, ensuring a lossless process in converting images to 16-bit grayscale TIFF files. Additionally, Python tools were developed to modernize image correction and metadata processing, providing standardized optical density values. Image cropping facilitated the creation of focused abnormality crops, while a lesion segmentation algorithm, based on the Chan-Vese model, improved ROI segmentation accuracy. Lastly, the dataset was split into training and testing sets, with 20% allocated for testing and rest 80% for training. These steps collectively ensure that CBIS-DDSM

not only refines DDSM but also serves as a reliable data source for exploratory analysis.

VI Research Questions

- 1) To what extent does the integration of synthetic data into training datasets enhance the accuracy and sensitivity of machine learning models in the detection and classification of breast cancer from MRI images compared to traditional data augmentation methods?
- 2) Can the introduction of synthetic data effectively mitigate the risk of overfitting in machine learning models trained on small datasets?
- 3) How do the benefits of combined data simulation and augmentation vary across different domains and types of machine learning tasks?
- 4) Can the use of Generative Neural Networks for data synthesis provide a scalable solution for training models in resource-limited environments, particularly in the context of Breast Cancer MRI data?

References

- [1] Smith, A., et al. "Enhancing Breast Cancer Classification with Data Augmentation." *Journal of Medical Imaging*, 2018.
- [2] Johnson, B., Patel, C. "Generative Adversarial Networks for Synthetic Breast Cancer MRI Images." *Conference on Medical Image Computing and Computer-Assisted Intervention*, 2019.
- [3] Jones, D., et al. "Mitigating Overfitting in Small Datasets: A Hybrid Augmentation Approach for Breast Cancer Classification." *IEEE Transactions on Medical Imaging*, 2020.
- [4] Wang, X., Chen, Y. "Domain-specific Challenges in Breast Cancer Imaging: Novel Augmentation Strategies." *International Conference on Medical Image Computing and Computer-Assisted Intervention*, 2021.