

International Telecommunication Union

ITU-T FG-AI4H Deliverable

TELECOMMUNICATION
STANDARDIZATION SECTOR
OF ITU

16 March 2023

PRE-PUBLISHED VERSION

DEL06

AI training best practices specification

ITU-T

Summary

Machine learning models for AI in Health are deployed in high-impact tasks. As a result, it is important to follow best practices for training and documentation so as to achieve maximum performance and transparency. The first part of this document provides a review of best practices for proper AI model training. The second part of this document provides guidelines for model reporting.

Keywords

AI training; best practices

Change Log

This document contains Version 1 of the Deliverable DEL06 on "*AI training best practices specification*" [approved at the ITU-T Focus Group on AI for Health (FG-AI4H) meeting held in 16 March 2023| approved on 16 March 2023 via the online approval process for the ITU-T Focus Group on AI for Health (FG-AI4H)].

Editor:	Sim Xinming AI Singapore	Email: xinming@aisingapore.org
	Stefan Winkler National University of Singapore	Email: winkler@nus.edu.sg

CONTENTS

	Page
1 Scope.....	1
2 References.....	1
3 Terms and definitions	1
3.1 Terms defined elsewhere	1
4 Abbreviations.....	1
5 Conventions	1
6 AI pipeline: a brief description	2
7 Best practices for data pre-processing	2
7.1 Feature engineering	2
7.2 Treatment of missing data	3
7.3 Data augmentation	4
8 Best practices for model training	4
8.1 Model architecture	4
8.2 Model ensembles	5
8.3 Model hyperparameters	5
8.4 Model validation	6
8.5 Model generalisability	8
8.6 Other important considerations.....	9
8.7 Model reporting	9
References	Error! Bookmark not defined.

List of Figures

	Page
Figure 1: A brief description of an AI training pipeline	2

ITU-T FG-AI4H Deliverable DEL06

AI training best practices specification

1 Scope

This deliverable provides a general introduction and background to model training across a typical AI pipeline. It explains basic key concepts, considerations, practices, and limitations on the different aspects of model training in the healthcare domain and points readers to cited sources or studies where advanced materials can be found.

2 References

[FG-AI4H DEL0.1] ITU-T FG-AI4H DEL0.1 (2022), *Common unified terms in artificial intelligence for health*.
https://www.itu.int/dms_pub/itu-t/opb/fg/T-FG-AI4H-2022-1-PDF-E.pdf

Additional references are found in the Bibliography clause of this document.

3 Terms and definitions

3.1 Terms defined elsewhere

The terms defined in [FG-AI4H DEL0.1] are applicable to this document.

3.2 Terms defined here

This document defines the following term:

3.2.1 Alarm fatigue: Also known as alert fatigue, is the sensory overload and desensitization that occurs when one is exposed to a large number of alerts and alarms in the health care context, leading to missed alarms or delayed response. (Adapted from Wikipedia and [34].)

4 Abbreviations

AI	Artificial intelligence
AutoML	Automated machine learning
EHR	Electronic health record
LOCF	Last observation carried forward
LR	Learning rate
MAR	Missing at random
MCAR	Missing completely at random
NMAR	Not missing at random
OHDSI	Observational Health Data Sciences and Informatics
OMOP	Observational Medical Outcomes Partnership
RNN	Recurrent Neural Networks

5 Conventions

This document does not use any particular conventions.

6 AI pipeline: a brief description

Given the complexities of diseases and their interaction with the human body, AI for health requires a high degree of customisation. In this document, we broadly cover the important considerations throughout an AI pipeline highlighted in Figure 1.

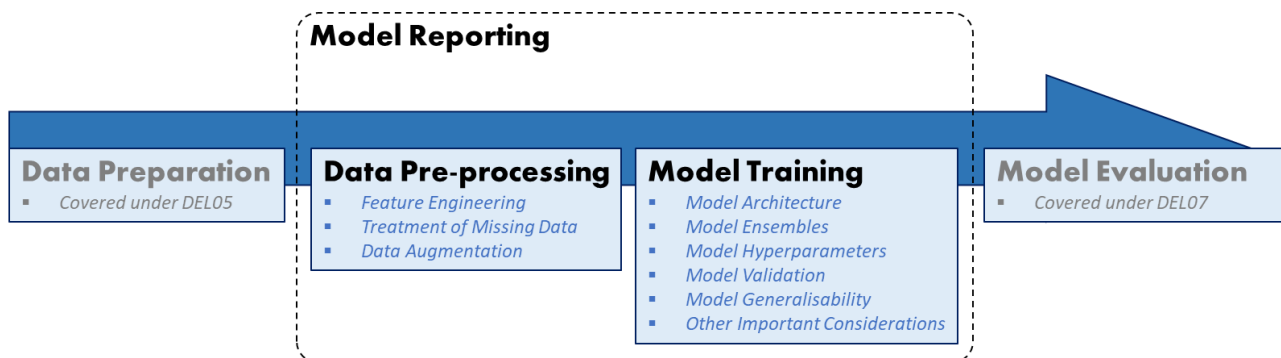


Figure 1: A brief description of an AI training pipeline

7 Best practices for data pre-processing

Data pre-processing is an important aspect in the AI model training pipeline. Healthcare data exist in many forms, such as physicians' notes, patient details, medication and medical images, which are typically consolidated and referred to as Electronic Health Record (EHR). Within EHR data, the number of domains (e.g., patient demographic, lifestyle, and health-related variable groups) and their individual sub-domains (e.g., occupation, smoking habits, exercise frequency and intensity, family history etc.) predictor variables, or features, present are large and complex.

However, not all features are equally significant to the problem statement, certain specific biomarkers hold a greater weightage, or are more clinically meaningful, in determining the outcome of a medical examination, e.g., specific cancer detection. It is imperative to scrutinize datasets, not just through the lens of AI or data scientists, but equally important also through the lens of medical domain experts in the fields around the problem statement. However, employing traditional clinical approaches and models to peruse all available dataset features would be tedious, time-consuming, and resources intensive. How then should the data be explored meticulously to identify and select data of significance?

7.1 Feature engineering

Before embarking on feature engineering, the problem statement must be clearly articulated to understand the desired outcomes. This highlights the necessary features, if present within the dataset, that contribute significantly to the desired outcomes, and elucidates other possible features that could be generated to complement them. As not all features are necessary, careful feature selection, reduction and generation is critical in building an accurate AI model, while saturating the AI model with irrelevant or inconsequential data degrades its performance or exacerbate resource costs. As AI applications proliferate within the healthcare space, we see a shift of approaches that guide feature selection or reduction.

- Traditional approaches. In traditional prediction models, predictor variables are typically predefined [1], such as blood pressure or cholesterol in the Framingham Heart Study [2]. In most cases, these features are selected based on expert judgement or knowledge. However, the depth of patients' data is not completely utilized, which may explain poorer model performances when compared to newer, automated approaches.
- Semi-automatic approaches. These approaches may consist of two components in varying degree; one which relies on automatic data-driven techniques to select and extract clinically meaningful features from large datasets, while complementing existing pre-defined features

selected manually based on expert judgement or established clinical studies and industry standards. This usually results in an AI model that outperforms traditional approaches [3].

- Deep learning approaches. Deep learning allows the AI model to execute feature engineering without being explicitly programmed, removing the need to pre-define features. Deep learning AI models are also able to handle large, messy, or incomplete datasets sufficiently well. In a study by Google AI on deep learning using EHR data [4], the utilisation of the full spectrum of EHR data using deep learning techniques has been shown to outperform traditional clinical models in prediction tasks such as in-hospital mortality, 30-day unplanned readmission, prolonged length of stay and final discharge diagnoses.

Feature generation is often done in tandem with feature selection. Additional features derived from the dataset is sometimes required to help the AI model learn better. Determining what additional features should or could be generated depends greatly on the problem statement and its relevant medical domain knowledge. For example, when determining the effectiveness of a particular treatment course and recommending treatment options to the healthcare professional for patient management, it is necessary to generate and provide AI models with data on historical information, also lag features, such as patient health markers over a time-period. It is therefore useful to always keep in mind a guiding framework on approaching feature extraction, generation, and aggregation [5] throughout the AI model training pipeline.

7.2 Treatment of missing data

It is unrealistic to assume that EHR data can be readily consumed for model training upon retrieval. It is common to expect missing values within healthcare records. How these missing data (e.g., missing true positive alarms) is managed, rationalized, and imputed could drastically impact the final model performance (e.g., resulting in false positives or false negatives).

Missing data could be in statistical terms classified into three main types: (1) missing completely at random (MCAR); (2) missing at random (MAR); and (3) not missing at random (NMAR). The easiest way to handle missing data is to simply ignore and remove either the value or the entire case itself. When doing so an underlying assumption that the missing data is independent of both the observed and unobserved data is made, which may not be the case. Several extensive methods used to address missing data in various healthcare context can be found in open literature [6], and are summarized below:

- Last observation carried forward (LOCF). LOCF carries forward observation from the last observed time point to fill the missing data at the endpoints.
- Complete-case analysis and available-case analysis. Complete-case analysis uses only data from patients with a complete record of all visits and ignores all patients with any missing data. Available-case analysis on the other hand uses all available data.
- Mean imputation, hot-deck imputation, regression imputation. Mean imputation uses the mean across all data at the time point to fill in the missing data. Hot-deck imputation fills the missing data with observed responses from another randomly selected but similar data point. Regression imputation uses a regression model to predict and fill the missing value.
- Mixed effects models and generalized estimating equations. For mixed effects models, a statistical distribution (e.g., gaussian distribution) is used to account for the missing data. On the contrary, an estimating function could also be used to predict the missing data.
- Inference for NMAR data. There are two methods that can be used to infer for NMAR data: pattern mixture models and selection models. Pattern mixture models [7] express the joint distribution of the responses and missing indicators using the distribution of all possible missingness patterns and the distribution of the responses given a specific missingness pattern. Selection models [8] express the joint distribution of the responses and the missing

indicators using the opposite decomposition; the distribution of the responses and the distribution of the missing indicators given the responses.

7.3 Data augmentation

Data augmentation is another important component of data pre-processing and has slight overlaps with feature generation. It is often used for computer vision applications when image data is limited or unbalanced (e.g., proportion of a certain data class is much higher or lower than the others), such as medical imaging and diagnosis applications. Data augmentation generates additional training data to complement existing datasets, providing the training process with more information to learn from, and boosting the model's accuracy, training stability or generalisability.

The types of data augmentation techniques that could be deployed depends on specific use-cases. Below are some basic techniques frequently deployed for image-based AI models, as well as references to more advanced techniques [9].

- Geometric Transformations. Techniques such as image cropping, zooming, rotating, translating are commonly used to create synthetic images to complement existing datasets. This helps expand dataset size while preserving its authenticity and may help to reduce positional biases during model training.
- Colour Space Transformations. Poor image lighting or colour balance may affect the model's ability to identify features during model training, which typically happens for medical imaging applications. Techniques such as equalising or normalising colour histograms often helps to improve model performance. However, caution must be exercised to ensure any colour manipulations do not affect key definitions within the image-data that characterises these features.

Data augmentation is less frequently deployed for text-based AI models compared to image-based ones due to a few factors. First, it is much easier to acquire large text-based data as compared to image-based data, hence less likely needing to overcome challenges such as limited data or class imbalance. Second, text-based data is much more complex due to the intricacies of the language medium. The resulting sentence sentiment, position and nuance could be drastically altered if one word is changed incorrectly. Yet, the benefits of data augmentation seen for image-based models has been recently demonstrated if they are used correctly for text-based models. Recent studies include Easy Data Augmentation (EDA) techniques [10] that could be deployed for text-based deep learning models, such as synonym replacement, random insertion, random swap, and random deletion.

8 Best practices for model training

Model training is not a linear process. Many iterations are required to properly train and optimise before the model sufficiently performs its intended functions. The international research community has explored countless aspects of model training over the past decades, but what is important to note is that the benefits of these developed techniques, tools, takeaways, conclusions from various studies are very dependent on the use-cases and settings, context, or environment they are applied to. As such, this section will only broadly focus on common basic best practices that help to bring a potential user up to speed on how to train and deploy an AI model.

8.1 Model architecture

The wide spectrum of AI models available is daunting, and the complexity, number and types of models will only proliferate further over time. It would be laborious for users with little background in AI to even begin to understand which types of models to use for a particular use case, much less evaluate their performance.

For experimentation purposes, it is advisable to start with AI models readily available and updated on open-source libraries (e.g., ResNet-152, GoogLe-Net19, BERT etc.) that have been trained on

similar datasets to achieve similar objectives, and conduct transfer-learning techniques to adapt them to the defined use-case. A survey of open-sourced deep learning architectures for various use cases provides some basic evaluations of each architectures, their strengths, and limitations [11]. Additionally, for more advanced models, there have been recent works summarizing different neural network architectures and their limitations for the healthcare domain based on their general use cases for bioinformatics, medical informatics, medical imaging, and public health [12].

While most of the document revolves around the basic application of AI, it is important to consider throughout the AI pipeline on the scalability implications due to factors such as model architecture. Where scalability becomes an issue, it may become necessary to adopt advanced techniques that could help to reduce model complexity while preserving model fidelity [13].

8.2 Model ensembles

It is a common misconception that AI applications are just one singular AI model working behind the scenes. Experimentation during the model training pipeline will show that some models are inherently better at certain tasks than others. Each model has its own strengths and weaknesses. It is difficult to optimise a model that is “good-at-everything”, which is usually limited by either an unfeasibly large amount of data required for it to learn from or the unrealistic amount of time needed to train it.

As such, the principle of model ensembles is to build multiple, smaller, implementable models and combining them together based on a set of decision-making or learning strategies. Multi-model ensembles have been demonstrated to achieve equal, if not better, performance than individual models for the detection of cancer [14], heart diseases [15] and diabetes prediction [16]. Here, we summarise basic decision-making techniques for model ensembles, as well as references to more advanced techniques on ensemble learning [17].

- Maximum or Majority Vote. Maximum or majority vote is employed by counting the prediction of each model. The prediction with the highest or majority vote will be used as the eventual ensemble prediction.
- Average or Weighted Vote. Average or weighted vote is employed by summing each models’ output probabilities and averaging them either uniformly (average vote) or based on a defined weightage (weighted vote). The output class with the highest probability will be used as the eventual ensemble prediction.

The type of ensemble methods to be used depends on factors such as model learning characteristics, domain knowledge etc. A study comparing the effectiveness of various ensemble methods and summarising the strengths and weaknesses of each type can be found here [18].

8.3 Model hyperparameters

Hyperparameter are a set of parameters that cannot be learnt by the model, and whose values are specified before model training begins to control the training process, such as learning rates, dropout rates, batch sizes, epochs etc. Control over these hyperparameters throughout the model training process is crucial to ensuring model accuracy, generalisability, and training stability (convergence). For e.g., high values of learning rate at the start of the model process may lead to divergence, while low values of learning rates increase the resource cost of model training. As such, it is common to see learning rate annealing techniques being used for better control over the training process.

Optimization of these hyperparameters ensures that the resultant model optimally minimises the predefined loss function for a given set of data. For professionals with little or no experience in machine learning, the process of optimizing these hyperparameters may be challenging. Furthermore, it becomes increasingly resource-intensive to do hyperparameter optimisation as the datasets increase in volume and complexity. Here we discuss some common manual and automated techniques used for hyperparameter optimization which could be used.

- Grid Search. Grid Search does an exhaustive search of a manually specified subset of the hyperparameter space. This is one of the simplest techniques and is easy to implement. However, it spends substantial time evaluating hyperparameter combinations which are unpromising and is therefore less efficient than Random Search.
- Random Search. Like the Grid Search, Random Search evaluates random combinations of hyperparameters within the search space and defined number of iterations. As not all combinations are evaluated over the search space, this makes hyperparameter optimization computationally less demanding when optimizing for higher dimensions of hyperparameters as compared to Grid Search.
- Bayesian Optimization. Bayesian optimization is a commonly adopted optimisation framework for many automated machine learning (AutoML) system [19]. It is an approach to optimizing objective functions that take a long time (minutes or hours) to evaluate. It is best suited for optimization over continuous domains of less than 20 dimensions and tolerates stochastic noise in function evaluations. It builds a surrogate for the objective and quantifies the uncertainty in that surrogate using a Bayesian machine learning technique, Gaussian process regression, and then uses an acquisition function defined from this surrogate to decide where to sample [20]. The difference between Bayesian Optimization methods from Grid Search or Random Search is that it uses past evaluation results to determine the subsequent values to evaluate. In essence this makes the optimization process more efficient as it limits iterations with poor hyperparameter combinations while focusing on promising hyperparameter combinations obtained from past results.

8.4 Model validation

Validating the AI model is a crucial step within the model training pipeline and occurs as the model trains. In this section, we are interested in knowing how well the model is performing to its intended objective based on dataset that was not yet seen by the model. It gives us an indication of the model's performance (how accurate are these decisions based on the input dataset?) and generalisability (how accurate are these decisions based on another unseen dataset?). To do so, we need to look out for instances of overfitting or underfitting within the AI model.

Overfitting occurs when the AI model or algorithm fits the data too well and unnecessarily captures the noise present within which affects the performance of the model when generalising to data beyond that used in its training. Underfitting, on the contrary, is when the AI model or algorithm cannot capture the trends or features of the data well enough. Both overfitting and underfitting results in poor model performance and could happen due to factors such as model parameters/weights, algorithm, and model selection.

- Selection of model or algorithms. It is impossible to represent a sine function well with a linear equation. Similarly, ensuring that the complexity of the model is befitting of the dataset affects how good the model outputs will be. In general, models or algorithms which are overly restrictive limits the learning of the dataset, and is likely to lead to an underfit, while models or algorithms which are overly complex captures unnecessary randomness within the training data and is likely to produce an overfit.
- Recognizing the warning signs. Knowing when the data is overfitted or underfitted is crucial in addressing the issue. Often, an overfit is typically the case compared to an underfit. There are two quick ways to identify overfit-underfit:
 - Bias-variance trade-off is one of the ways to know when the best fit can be achieved for a given setting. Bias is an error due to erroneous assumptions and is an indication of a possible underfit. Conversely, variance refers to the model or algorithm's sensitivity to perturbations within the datasets. High variance generally occurs when the model or algorithm is overly complex and captures the randomness with the sample, giving an

indication of a possible overfit. To ensure a well- performing model, one needs to find right balance between bias and variance which minimizes the total error function.

- Training and validation performance plots over the training cycles are another common indicator to look at to determine whether an AI model is overfitting its data. The performance (accuracy or loss) plots indicate how well the model is performing at that given training cycle. One clear indication of overfitting is when the validation performance starts to deteriorate (accuracy decreases or loss increases) and deviate from its converged training performance over the number of training cycles.
- Cross-validation to prevent overfits. While there are other methods which could detect and prevent overfitting such as back-testing or regularization, one of the more popular methods is cross validation. Though there are many different variations of cross validation methods such as k-fold or stratified k-fold cross validation, the underlying principles are the same, which is to partition the data into (1) training subset and (2) validation subset. At any iteration, the model only trains and tunes on the training subset. The validation on the validation subset gives an indication of the performance on unseen (test) data. The final evaluation of the model on the completely unseen test dataset gives an unbiased assessment of its performance. If the model outperforms on the training data subset compared to the test set, it is likely to be an indicator that the model is producing an overfit.
- Number of model parameters/weights. The number of parameters affects the resulting performance of the model. Having too few leads to underfitting, while conversely having too many may lead to an overfit. It is recommended for parameter optimization to be carried out to determine what is the optimal number to be used to train the model. Underfitting can be resolved by increasing the number of parameters or using weaker regularization during model training.

Overfitting is often encountered during model training, compared to underfitting. Underfitting is generally an easier issue to tackle with simple solutions such as increasing model complexity, increasing model training over more epochs. However, tackling overfitting requires a greater depth of understanding in both the dataset and the model algorithms. Below are a few fundamental strategies that could be employed, as well as references to more advanced techniques [21].

- Simplifying model complexity. Neural networks can sometimes be prone to overfitting. While there is no single hard rule or guideline to follow, in instances of an overfit, one should aim to reduce the complexity of the network, such as reducing the number of hidden layers or reducing the number of neurons.
- Employing regularization techniques. There are many different types of regularization techniques developed for different types of AI models, but their main objective is to limit the capacity of the AI by penalizing the model training loss functions (e.g., L1 or L2 regularization) or by reducing the complexity of the algorithm through dropouts [22]. Recent work around some common and novel regularization techniques for computer vision applications and their effectiveness can also be found here [23].
- Early stopping. As discussed in the previous section, one of the common indicators of overfitting is when the validation performance deteriorates and deviates from the converged training performance. As such, early stopping of model training (or rolling back of model weights to a saved point before overfitting occurs) is a strategy that is always used due to its simplicity.
- Data augmentation. Highly unbalanced datasets with disproportionate distribution of classes (e.g., in applications where patients at high risk of a disease belongs to a minority class) should be carefully managed to prevent overfitting. One common best practice is to increase the training samples through data augmentation.

8.5 Model generalisability

So far, we have discussed the various best practices for model training. As model training is a highly iterative process, we need to identify useful and significant metrics to evaluate its performance besides model accuracy, i.e., accuracy of the model predictions with the ground truth. In the healthcare context, due to complexities across populations (e.g., some population groups are more susceptible to certain types of diseases), disease progressions (e.g., differing disease progression across various geographic regions), healthcare infrastructure among other factors, it is also critical to test the AI model on its generalisability across various factors before deployment and address them accordingly in the training process.

The generalisability of AI models has always been a key challenge in the community, and many are far from achieving reliable generalisability. In the healthcare context, AI models are relatively constrained in their deployment to populations and settings like those they were trained on. For example, diseases may manifest itself slightly different across ethnicity. What happens if the AI models were trained on a particular ethnic distribution and tested on a different one? While there has been research showing some success in developing limited generalisable AI models across ethnic distributions [24], this remains a key concern for healthcare professionals in ensuring that the diagnosis made by the AI model is as accurate as it is expected to be.

- Generalisability across different populations. Several measures have been proposed to help overcome the issue of AI generalisability across different populations [25]:
 - Site-specific training to adapt an existing system for a new population, particularly for complex tasks like EHR predictions. Methods to detect out-of-distribution inputs [26][27] and provide a reliable measure of model confidence to prevent clinical decisions being made on inaccurate model outputs. For medical image classification, this problem may overcome by the curation of large, heterogenous, multi-centre datasets. However, the limitation to this method is that re-training becomes computationally demanding and expensive should the number of out-of-distribution cases far exceed a certain threshold.
- Generalisability across different settings. Healthcare institutions around the world vary in their protocols, hardware such as medical equipment and medical equipment software. This results in vastly different healthcare data generated by each unique healthcare institution, which makes it a nightmare to generalise AI system from one institution to another. While there have been efforts by organisations such as Observational Health Data Sciences and Informatics (OHDSI) to standardize a common data, model called Observational Medical Outcomes Partnership (OMOP) for electronic healthcare recording, this remains a barrier to the interoperability of AI models across institutional and international boundaries.
 - One possible way of overcoming this issue lies with the design of the AI model and training. Splitting a disease diagnosis step into segmentation and measurement may enable easier generalization to new imaging hardware by retraining only the segmentation model, which is more data- efficient [28].
- Generalisability across time. Besides the issue of adapting an AI system to new populations, one must also consider the temporality and everchanging nature of diseases. Diseases are always progressing and changing over time in a nondeterministic way [29]. However, many existing deep learning models, including those already proposed in the medical domain, assume static vector-based inputs, which cannot handle the time factor in a natural way.
 - To address this, it is recommended to design deep learning approaches that can handle temporal health care data. In this aspect, Recurrent Neural Networks (RNN) is a good choice to utilize when healthcare data is sequentially ordered, as they are well designed to handle temporal dependencies [30].

8.6 Other important considerations

Beyond that however, there are other important considerations that should be considered as well throughout the model training pipeline.

- Trust in AI model and systems. Deep learning has been extensively explored in the medical imaging space, however not many AI systems have been deployed in the real-world context. Despite showing high performance during model evaluations, a lack of trust by clinicians towards these systems still exist. More recent works documenting proposals to address this include using techniques such as incorporating uncertainty modelling into the AI system [31] so that clinicians are better informed of how certain the AI system is when deriving a certain decision outcome, as well as frameworks to provide patient-level interpretation derived from AI model features to help with patient management, as well as feature-level interpretation derived from aspects within the datasets, which could help with medical research [32].
- Model validation source. Not all data are created equal. How the model was validated may be a factor to consider in the healthcare context. Healthcare data from prospective studies generally have fewer biases and are more favourable for model validation than from retrospective studies. This however remains a challenge, as the number of established prospective studies are limited. Models that are trained and validated from clinical studies spanning multiple centres are also more likely to be generalisable, allowing the model to be adopted for a different population group significant model retraining effort.
- Workflow integration and implementation. Developing a superior but overly complex AI model is worthless if it does not fit well within the clinician's workflow and eventually ends up being side-lined, i.e., requiring additional user interface, software, or hardware to integrate the AI solution within the hospital. If the AI model excessively prompts the clinician with alarms, it may lead to "alarm fatigue", leading to these alerts being ignored or switched off.
- Clinical applicability and impact. AI models are typically assessed against a baseline performance. However, for the healthcare context, such systems should be evaluated against a human clinician to determine if it does indeed result in an impact to the level of care received by a patient. Even so, a well performing model may not necessarily translate to an overall higher clinical impact. The system must be assessed prospectively in practice to determine the level of impact it can deliver to clinicians and patients.

While there are many other important aspects that need to be considered as well, such as model privacy, security, and especially ethical considerations, these will not be discussed in detail here.

8.7 Model reporting

Now that the pipeline for AI model training has been described in clause 7, responsible development protocols dictate that proper recording and accounting framework be in place to ensure that the AI model is used within its boundaries and limitations. This ensures that users fully appreciate the AI models' purpose, performance, and particularly, limitations so that they are appropriately deployed.

Proper accounting and reporting procedures are critical for the deployment of AI models in healthcare, where model failures or misuse can lead to serious ramifications. As best practice, this should be done to encompass important information throughout the AI pipeline, rather than covering only model training.

This section describes a framework based on the "Model Cards" [33] proposed recently to encourage transparent reporting about a trained machine learning model. The purpose of this reporting framework is to enable those considering the use of a specific trained model in a particular context to better understand the systematic impacts of the model before deploying it. A standardized reporting format not only makes it easier for different stakeholders to assess and compare candidate models, but also encourages forward-looking model analysis in the development phase. An extract from the paper is summarized below:

- Model details. Basic information about the model, e.g., persons and/or organization developing model, funding, model development date, model version (and change log), model type, model outputs, sources and references, citation details, license details, and feedback channels.
- User case and users. Use cases that were envisioned during development, e.g., primary intended uses, primary intended users as well as out-of-scope uses.
- Factors. Factors could include demographic or phenotypic groups, environmental conditions, technical attributes, or others.
- Metrics. Metrics should be chosen to reflect potential real-world impacts of the model, e.g., model performance measures, decision thresholds and approaches towards uncertainty and variation.
- Training and evaluation data. Details on the dataset(s) used for the quantitative analyses, e.g., datasets, motivations, and pre-processing.
- Quantitative analyses. This provides the analyses and results of the model evaluation according to the metrics, and its corresponding confidence intervals, e.g., unitary analyses and intersectional analyses.
- Ethical considerations. Provide information on ethical considerations that went into model development, surfacing ethical challenges and solutions to stakeholders on issues such as data, risk management, risks and harms, use-cases, and the review processes.

Bibliography

- [1] Goldstein, B. A., Navar, A. M., Pencina, M. J., & Ioannidis, J. (2017). Opportunities and challenges in developing risk prediction models with electronic health records data: a systematic review. *Journal of the American Medical Informatics Association*, 24(1), 198-208.
- [2] Wilson, P. W., D'Agostino, R. B., Levy, D., Belanger, A. M., Silbershatz, H., & Kannel, W. B. (1998). Prediction of coronary heart disease using risk factor categories. *Circulation*, 97(18), 1837-1847.
- [3] Luo, G. (2019). A roadmap for semi-automatically extracting predictive and clinically meaningful temporal features from medical data for predictive modeling. *Global transitions*, 1, 61-82.
- [4] Rajkomar, A., Oren, E., Chen, K., Dai, A. M., Hajaj, N., Hardt, M., ... & Dean, J. (2018). Scalable and accurate deep learning with electronic health records. *NPJ Digital Medicine*, 1(1), 18.
- [5] Tran, T., Luo, W., Phung, D., Gupta, S., Rana, S., Kennedy, R. L., ... & Venkatesh, S. (2014). A framework for feature extraction from hospital medical data with applications in risk prediction. *BMC bioinformatics*, 15(1), 425.
- [6] Wong, W. K., Boscardin, W. J., Postlethwaite, A. E., & Furst, D. E. (2011). Handling missing data issues in clinical trials for rheumatic diseases. *Contemporary clinical trials*, 32(1), 1-9.
- [7] Little, R. J. (1993). Pattern-mixture models for multivariate incomplete data. *Journal of the American Statistical Association*, 88(421), 125-134.
- [8] Diggle, P., & Kenward, M. G. (1994). Informative drop-out in longitudinal data analysis. *Journal of the Royal Statistical Society: Series C (Applied Statistics)*, 43(1), 49-73.
- [9] Shorten, C., & Khoshgoftaar, T. M. (2019). A survey on image data augmentation for deep learning. *Journal of Big Data*, 6(1), 60.
- [10] Wei, J., & Zou, K. (2019). Eda: Easy data augmentation techniques for boosting performance on text classification tasks. *arXiv preprint arXiv:1901.11196*.
- [11] Alom, M. Z., Taha, T. M., Yakopcic, C., Westberg, S., Sidike, P., Nasrin, M. S., ... & Asari, V. K. (2019). A state-of-the-art survey on deep learning theory and architectures. *Electronics*, 8(3), 292.
- [12] Zion, I., Ozuomba, S., & Asuquo, P. (2020, March). An Overview of Neural Network Architectures for Healthcare. In *2020 International Conference in Mathematics, Computer Engineering and Computer Science (ICMCECS)* (pp. 1-8). IEEE.
- [13] Xu, D., Lee, M. L., & Hsu, W. (2019). Propagation Mechanism for Deep and Wide Neural Networks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition* (pp. 9220-9228).
- [14] Xiao, Y., Wu, J., Lin, Z., & Zhao, X. (2018). A deep learning-based multi-model ensemble method for cancer prediction. *Computer methods and programs in biomedicine*, 153, 1-9.
- [15] Mienye, I. D., Sun, Y., & Wang, Z. (2020). An improved ensemble learning approach for the prediction of heart disease risk. *Informatics in Medicine Unlocked*, 20, 100402.

- [16] Akula, R., Nguyen, N., & Garibay, I. (2019, April). Supervised Machine Learning based Ensemble Model for Accurate Prediction of Type 2 Diabetes. In 2019 SoutheastCon (pp. 1-8). IEEE.
- [17] Dong, X., Yu, Z., Cao, W., Shi, Y., & Ma, Q. (2020). A survey on ensemble learning. *Frontiers of Computer Science*, 1-18.
- [18] Ju, C., Bibaut, A., & van der Laan, M. (2018). The relative performance of ensemble methods with deep convolutional neural networks for image classification. *Journal of Applied Statistics*, 45(15), 2800-2818.
- [19] Waring, J., Lindvall, C., & Umeton, R. (2020). Automated machine learning: Review of the state-of-the-art and opportunities for healthcare. *Artificial Intelligence in Medicine*, 101822.
- [20] Frazier, P. I. (2018). A tutorial on bayesian optimization. *arXiv preprint arXiv:1807.02811*.
- [21] Moradi, R., Berangi, R., & Minaei, B. (2019). A survey of regularization strategies for deep models. *Artificial Intelligence Review*, 1-40.
- [22] Srivastava, N., Hinton, G., Krizhevsky, A., Sutskever, I., & Salakhutdinov, R. (2014). Dropout: a simple way to prevent neural networks from overfitting. *The journal of machine learning research*, 15(1), 1929-1958.
- [23] Xu, D., Lee, M. L., & Hsu, W. (2019, September). Patch-Level Regularizer for Convolutional Neural Network. In 2019 IEEE International Conference on Image Processing (ICIP) (pp. 3232-3236). IEEE.
- [24] Bellemo, V., Lim, Z. W., Lim, G., Nguyen, Q. D., Xie, Y., Yip, M. Y., ... & Ting, D. S. (2019). Artificial intelligence using deep learning to screen for referable and vision-threatening diabetic retinopathy in Africa: a clinical validation study. *The Lancet Digital Health*, 1(1), e35-e44.
- [25] Kelly, C. J., Karthikesalingam, A., Suleyman, M., Corrado, G., & King, D. (2019). Key challenges for delivering clinical impact with artificial intelligence. *BMC medicine*, 17(1), 195.
- [26] Liang, S., Li, Y., & Srikant, R. (2017). Principled detection of out-of-distribution examples in neural networks. *arXiv preprint arXiv:1706.02690*, 655-662.
- [27] Hendrycks, D., & Gimpel, K. (2016). A baseline for detecting misclassified and out-of-distribution examples in neural networks. *arXiv preprint arXiv:1610.02136*.
- [28] De Fauw, J., Ledsam, J. R., Romera-Paredes, B., Nikolov, S., Tomasev, N., Blackwell, S., ... & Ronneberger, O. (2018). Clinically applicable deep learning for diagnosis and referral in retinal disease. *Nature medicine*, 24(9), 1342-1350.
- [29] Miotto, R., Wang, F., Wang, S., Jiang, X., & Dudley, J. T. (2018). Deep learning for healthcare: review, opportunities and challenges. *Briefings in bioinformatics*, 19(6), 1236-1246.
- [30] Shickel, B., Tighe, P. J., Bihorac, A., & Rashidi, P. (2017). Deep EHR: a survey of recent advances in deep learning techniques for electronic health record (EHR) analysis. *IEEE journal of biomedical and health informatics*, 22(5), 1589-1604.
- [31] Lim, Z. W., Lee, M. L., Hsu, W., & Wong, T. Y. (2019, July). Building Trust in Deep Learning System towards automated disease detection. In *Proceedings of the AAAI Conference on Artificial Intelligence* (Vol. 33, pp. 9516-9521).

- [32] Zheng, K., Cai, S., Chua, H. R., Wang, W., Ngiam, K. Y., & Ooi, B. C. (2020, June). TRACER: A Framework for Facilitating Accurate and Interpretable Analytics for High Stakes Applications. In Proceedings of the 2020 ACM SIGMOD International Conference on Management of Data (pp. 1747-1763).
 - [33] Mitchell, M., Wu, S., Zaldivar, A., Barnes, P., Vasserman, L., Hutchinson, B., ... & Gebru, T. (2019, January). Model cards for model reporting. In Proceedings of the conference on fairness, accountability, and transparency (pp. 220-229).
 - [34] Woo M, Bacon O. Alarm Fatigue. In: Hall KK, Shoemaker-Hunt S, Hoffman L, et al. Making Healthcare Safer III: A Critical Analysis of Existing and Emerging Patient Safety Practices [Internet]. Rockville (MD): Agency for Healthcare Research and Quality (US); 2020 Mar. 13.
<https://www.ncbi.nlm.nih.gov/books/NBK555522/>
-