

Inducing natural invariance in deep segmentation pipelines for generalizable detection of colorectal polyps

Any short subtitle

Birk Sebastian Frostelid Torpmann-Hagen



Thesis submitted for the degree of
Master in Robotics and Intelligent Systems
60 credits

Department of Informatics
Faculty of mathematics and natural sciences

UNIVERSITY OF OSLO

Spring 2022

Inducing natural invariance in deep segmentation pipelines for generalizable detection of colorectal polyps

Any short subtitle

Birk Sebastian Frostelid Torpmann-Hagen

© 2022 Birk Sebastian Frostelid Torpmann-Hagen

Inducing natural invariance in deep segmentation pipelines for
generalizable detection of colorectal polyps

<http://www.duo.uio.no/>

Printed: Reprosentralen, University of Oslo

Abstract

Contents

1	Introduction	1
	Introduction	1
2	Background	3
2.1	Colorectal Polyps, Medical Imaging, and Deep Learning . .	3
2.2	Generalization failure in the Wild	4
2.3	Generalizability Theory	5
2.3.1	Generalizability and Expectations	5
2.3.2	Empirical Risk Minimization	5
2.3.3	Underspecification	6
2.3.4	Generalisation Failure Modes	7
2.3.5	Bayesian perspective of Generalization	7
3	Methodology	9
	Methodology	9
3.1	(algorithm name)	9
3.1.1	Model of natural variability	9
3.1.2	Geometric and pixel-wise transformations	10
3.2	Baselines	10
3.3	Datasets	10
3.4	Metrics and evaluation	10
4	Results	11
	Results	11
5	Discussion	13
	Discussion	13

List of Figures

List of Tables

Preface

Chapter 1

Introduction

Colorectal cancer is one of the leading causes of cancer related deaths, causing approximately 900 thousand deaths worldwide per year (cite). Early detection thereof is as a consequence of significant importance. Polyps are often an early warning-sign of developing tumor, and early detection thereof can as a result significantly reduce fatality rates. Polyps are, however, often missed during colonoscopies, owing to the significant variability in the shape and size of polyps, as well as the high degrees of similarity to surrounding tissue. Automatic segmentation of polyps via deep learning has the potential to significantly increase the likelihood of early detection and effective treatment.

However, clinical applications of deep learning are known to fail in deployment, despite exhibiting excellent performance during development. This is known as generalization failure, and is ubiquitous in the domain. While there has been a growing body of research dedicated to identifying and analyzing its root causes, there is still limited research into approaches to mitigating generalizability failure.

This thesis presents a novel approach to increasing generalizability, whereby the model is trained to not only minimize segmentation-loss, but also minimize the effects of the data being perturbed by an ensemble of transformations, including color-transformations, geometric transformations, additive noise, and adding extra polyps to the image using a GAN-inpainter. This endows the pipeline with the ability to more readily infer causally viable inductive biases by explicitly forcing the model to be robust to any combination of the aforementioned transformations.

Generalizability is then measured by evaluating several vanilla-pipelines consisting of several models on a number of separate datasets, which is then compared to the results of the modified pipeline. The results show that (...)

Chapter 2

Background

2.1 Colorectal Polyps, Medical Imaging, and Deep Learning

Polyps are small growths found in and around the inner lining of the large intestine. These polyps, also referred to as adenomas, can in time develop into cancerous tumours, or carcinomas, in a process known as the adenoma-carcinoma sequence [5]. Though the majority of polyps do not undergo this process, identifying polyps nonetheless constitutes an important step towards preventing colorectal cancer. Indeed, resection of these polyps has been shown to reduce the incidence of colorectal cancer by a significant margin [9].

Though colorectal cancer remains as one of the leading causes of cancer-related death worldwide (source), mortality rates have nonetheless declined in large part to increased use of screening colonoscopy, which in turn allows for the opportunity for preemptive treatment. Polyps are, however, by nature somewhat difficult to detect and are routinely missed by clinicians, with miss rates ranging upwards of 27% for diminutive polyps [4, 6].

Naturally, reducing this miss rate has the potential to further reduce incidence rates. There has been a significant body of work dedicated to for instance workflow optimization using both optical and mechanical approaches, for instance chromoendoscopy, wherein stains or dyes are applied at the time of endoscopy, or the use of alternative lighting such as narrow-band imaging, which increases the textural details that help distinguish polyps from their surrounding tissue.

These systems do, however, require more equipment, training and expertise to effectively employ. Thus, automatic polyp segmentation using deep learning and convolutional neural networks (CNNs) has been identified as another diminutive detection method. This requires minimal training time on the part of the clinician, no additional equipment, and has been shown to significantly increase detection rates when deployed in a clinical setting [1].

Naturally, these results are not unique to the detection of polyps. Indeed, medical imaging has in recent years proven to be one of the most

promising applications of artificial intelligence and deep learning, having the capacity to significantly improve both the accuracy and efficiency of detection, diagnosis, and treatment of a wide variety of diseases [8].

There are, however, still several hurdles to overcome; recent research has shown that even state of the art deep-learning pipelines are prone to generalization failure when deployed in practical settings, particularly when exposed to distributional shifts such as changes in demographics, imaging equipment, noise, and more despite exhibiting high performance on hold-out sets [2, 3]. There is no reason to expect that polyps are exempt from this problem, given how pervasive such shortcomings are in similar tasks.

Naturally, such systems are rendered practically useless should they fail to perform sufficiently outside of the very carefully controlled conditions upon which they are trained. Thus, for AI-assisted detection to be on any considerable merit, it has to infer causally reasonable patterns in the data that generalize well to other hospitals, demographics, imaging equipment, resolutions, and so on. Though a human would not find this type of generalization very difficult, deep-learning (and for that matter, other data-driven approaches) regularly seem to fail in this regard.

2.2 Generalization failure in the Wild

The medical domain is characterized by a number of key features that separate it from other areas where deep-learning typically excels. Training data is often scarce, the pathologies that constitute the classification targets are unevenly distributed and often exhibit high degrees of inter-class variability, and there can be a significant number of confounding variables.

For instance, a deep-learning based classifier which successfully detected pneumonia in X-ray scans across a number of hospitals with striking accuracy was determined to be basing its predictions not on any lesions or otherwise pathologically relevant features in the images, but rather on a hospital-specific metal token that was on every image, which it used in conjunction with learning the prevalence rate of pneumonia for the hospitals from which the data was collected. As a result, when deployed on data from hospitals that it had not seen during training, the system failed to generalize [10].

In another study, it was shown that a classifier intended to detect diabetic retinopathy exhibited significant variability in performance depending on the type of camera used. The same study also showed that the same type of performance variability could be found when detecting skin-conditions across demographics with differing skin tones [3].

Though there has been limited literature detailing the generalizability of pipelines trained to segment polyps, there is no reason to believe that it is somehow exempt from the same problems that impact the aforementioned tasks. To illustrate this, a brief meta analysis can be performed. Table There has been astonishingly little research into determining the generalizability of models designed for polyp-detection.

Model	IID	Mean OOD
DeepLabV3
U-Net
FPN
DDA-net
Divergent-net

As a result, a meta-analysis of a selection of such models is required. Several models, consisting of DDANet, DivergentNet, (...), were trained according to the hyperparameters provided in their respective papers and repositories. Table (...) shows their results when evaluated on separate datasets. Methodological details can be found in chapter (...).

Naturally, non-medical domains are in no way immune to generalisation failure. In fact, one could easily argue that the vast majority of deep-learning pipelines fail to generalize altogether, and instead merely infer some set of inductive biases that, although perhaps causally incorrect, perform sufficiently well for general use. Consider for example

2.3 Generalizability Theory

2.3.1 Generalizability and Expectations

When developing deep-learning based systems, there is by and large an expectation that it will

2.3.2 Empirical Risk Minimization

At the most fundamental level, the goal of machine learning is to learn a mapping between two spaces of objects X and Y . This mapping, namely the function $f : X \rightarrow Y$, maps some input object $x \in X$, an image for example, to a corresponding and application-relevant output object $y \in Y$, for instance a segmentation mask or a class label. It is worth noting, however, that f is not as much a function in the mathematical sense as much as it is an abstraction of whatever ground-truth relationship that the deep learning system is intended to capture, and consequently cannot typically be modelled explicitly. Instead, machine learning systems aim to find a representation of this mapping automatically by leveraging a training set $\{x_i, y_i\}_{0 \dots n}$ to find a sufficiently performant approximation of f . This is referred to as supervised learning, and the resulting approximation found using the training set is denoted by $h : X \rightarrow \hat{Y}$, and typically referred to as a hypothesis.

To find such an approximation, we assume that there exists a joint probability distribution over X and Y , namely $P(x, y)$, and that the training data $\{x_i, y_i\}_{0 \dots n}$ is drawn from this probability distribution such that the resulting sample distribution is independent and identically distributed (henceforth: iid) to $P(x, y)$. This is the so-called iid assumption. Note that by modelling the mapping as a joint probability distribution, one can model

uncertainty in the predictions by expressing the output as a conditional probability $P(y|x)$. In conjunction with a loss-function $L(h(x), y)$ which measures the discrepancy between the hypothesis and the ground truth, these assumptions allows us to quantify the expected performance of a given hypothesis:

$$R(h) = E[L(h(x), y)] = \int L(h(x), y) dP(x, y) \quad (2.1)$$

Using this framework, one can then find an iid-optimal hypothesis, often called a predictor, by finding the predictor h^* among a fixed class of functions (defined by network architecture) \mathcal{H} that minimizes risk:

$$h^* = \arg \min_{h \in \mathcal{H}} R(h) \quad (2.2)$$

Since $P(x, y)$ is not known, however, one cannot compute $R(h)$ explicitly. Instead, the expected risk has to be computed through empirical estimation, i.e by finding the arithmetic average of the risk associated with each prediction by the hypothesis over the training set:

$$R_{emp}(h) = \frac{1}{n} \sum_{i=1}^n L(h(x_i), y_i) \quad (2.3)$$

This risk can in turn be minimized with respect to the hypothesis class. This is called empirical risk minimization (ERM):

$$\hat{h} = \arg \min_{h \in \mathcal{H}} R_{emp}(h) \quad (2.4)$$

The central idea with this approach to machine learning is that the training data can be considered a finite iid sampling of the underlying distribution. As such, by the central limit theorem, the hold-out performance of the computed hypothesis will approach iid-optimal performance given a sufficient amount of training data and some sufficiently capable and regularized training procedure. This should in theory allow deep learning systems to be able to generalize, since the empirical risk in theory can approximate the true risk arbitrarily well given sufficient training data.

2.3.3 Underspecification

Naturally, however, real-world data is rarely neat enough for it to consistently abide by the iid assumption. Commonly encountered variation in real world data such as variable lighting conditions, class imbalance, image corruptions, noise, or other more subtle forms of distributional shift all result in structural misalignment of the training and deployment distributions (citation). Ideally, predictors should be robust to these sorts of changes, however evidently this is not guaranteed by ERM (citation). ERM simply guarantees an iid-optimal predictor. While the difference is subtle, it is worth reemphasizing: empirical risk minimization only generalizes to data which is more or less identically distributed to the training

data. Differently distributed or otherwise perturbed data, even that which is near imperceptible or at any rate inconsequentially different to the human eye, violates the iid assumption, and can as such not be expected to be classified correctly given a predictor trained via ERM.

To mitigate this, one could simply add more data to the pipeline through augmentation, or simply collecting more training data. This will lead to a better approximation of the true risk. This does not, however, solve the problem. The variability of the real world is not, unfortunately, easy to model merely through augmentations, and collecting sufficient data to cover every potential source of natural variability is infeasible, especially in medical domains. Consider for example a machine-learning pipeline wherein a model is trained to classify cows and camels. The dataset consists of cows, pictured in grass fields and pastures, and camels, pictured in deserts. To be generous, let us assume that we have sufficient quantities of data to ensure that the pipeline is perfectly invariant to the pose of the respective animals, to lighting conditions, geometric transforms, etc. One may then expect that the pipeline correctly learns to classify the two, and attains high accuracies, and indeed when evaluated on iid data, this would be entirely correct. However, what would then happen if one such predictor encountered a cow in the desert and a camel in a grass pasture? This constitutes a distributional shift, and as such we cannot expect reliable performance as detailed in ???. Naturally, the predictor may have learned just fine exactly what constitutes a cow and a camel, but it might just as easily learned to associate deserts with camels and pastures with cows. And from a data perspective, both are equally correct interpretations. The immediate response to this may be to simply add some pictures with more varied backgrounds, but this once again would only serve to make the pipeline more robust to backgrounds. It would not guarantee that the pipeline learns the right inductive biases. The predictor may then for example instead learn that cows typically are black and white and camels usually beige, and then fail when it encounters a brown cow. One could keep adding more and more data, but there is not really any way of knowing when the pipeline is well enough specified by the data such that it starts returning predictors with the desired inductive biases. There are in simpler terms several "correct" interpretations of what separates the classes from a purely data-based perspective, each with their own inductive biases. There are as a consequence not just one risk-minimizing predictor, but a whole family of them. This is referred to as underspecification [3].

2.3.4 Generalisation Failure Modes

2.3.5 Bayesian perspective of Generalization

Chapter 3

Methodology

As described in earlier sections, good generalizability can only be achieved if the pipeline can reliably produce predictors that infer the right inductive biases. Naturally, the set of correct inductive biases are not known, so any such pipeline instead has to learn to not infer the wrong inductive biases. To achieve this, a model of natural variance is constructed, which aims to encapsulate all the variability one might expect to see in the domain. This model can then be leveraged to force the pipeline to be robust to natural variance through contrastive learning. The central idea, then, is that it is more likely that the model learns to infer generalizable inductive biases as opposed to learning to simply be robust to all possible configurations of a large amount of transformations.

To evaluate this, several predictors are trained from several pipelines with and without the influence of (algorithm name). Their performance is then evaluated on both a stress-test, and two separate polyp datasets, namely Etis-larib and EndoCV2021).

3.1 (algorithm name)

3.1.1 Model of natural variability

In order to account for any natural variation one may expect to find in deployment, it is necessary to construct a model which can parameterize the variability that is encountered, in other words a model of natural variability (MNV). Naturally, there is no way of knowing the full extent thereof, but it may be sufficient to model some subset of the possible distributional shifts. This, naturally, requires some knowledge of the domain from which the dataset is collected. Similarly to how adding rotational augmentations is a bad idea for classification of hand-written numbers, certain transformations may or may not be suitable for use within a MNV.

In the case of polyp-segmentation, it is clear that it is necessary to account for variability in for instance lighting, polyp-size, polyp-shape, polyp-location, camera-quality, color-shifts, blurs, optical distortions, and affine transformations. Thus, a model is required that can (more or less)

parametrize this variability. Broadly speaking, these transformations can be categorized as follows:

- Pixel-wise variability, which affect only the image, i.e color-shifts, brightness shifts, contrast-shifts, lighting, blurs etc
- Geometric variability, which affect both the image and the segmentation mask by some parametrizable quantity, i.e affine transforms and distortions
- Manifold variability, which affects both the image and the segmentation mask depending on a learned model of the distribution, i.e the size, shape and location of polyps

Pixel-wise variability and geometric variability can be modeled fairly trivially through the use of the same transformations typically used for data-augmentation. Manifold-variability, however, is somewhat more difficult. Similar to how [7] employs cross-dataset style-transfer, it is necessary to find some way to model the distributional properties of the data, and then apply perturbations using the resulting model. Since both the size, shape, and position of polyps can be expected to vary, a model that can change all these factors is necessary. To this end, an in-painting model can be constructed. In particular, a GAN-inpainter.

Gan-based polyp inpainting

3.1.2 Geometric and pixel-wise transformations

3.2 Baselines

Several models were tested (...)

3.3 Datasets

3.4 Metrics and evaluation

Chapter 4

Results

Chapter 5

Discussion

y

asdfasdf

Bibliography

- [1] Ishita Barua et al. 'Artificial intelligence for polyp detection during colonoscopy: a systematic review and meta-analysis'. en. In: *Endoscopy* 53.3 (Mar. 2021), pp. 277–284.
- [2] Emma Beede et al. 'A Human-Centered Evaluation of a Deep Learning System Deployed in Clinics for the Detection of Diabetic Retinopathy'. In: *Proceedings of the 2020 CHI Conference on Human Factors in Computing Systems*. CHI '20. Honolulu, HI, USA: Association for Computing Machinery, 2020, 1–12. ISBN: 9781450367080. DOI: 10 . 1145 / 3313831 . 3376718. URL: <https://doi.org/10.1145/3313831.3376718>.
- [3] Alexander D'Amour et al. *Underspecification Presents Challenges for Credibility in Modern Machine Learning*. 2020. arXiv: 2011 . 03395 [cs.LG].
- [4] D Heresbach et al. 'Miss rate for colorectal neoplastic polyps: a prospective multicenter study of back-to-back video colonoscopies'. en. In: *Endoscopy* 40.4 (Apr. 2008), pp. 284–290.
- [5] A Leslie et al. 'The colorectal adenoma-carcinoma sequence'. en. In: *Br. J. Surg.* 89.7 (July 2002), pp. 845–860.
- [6] D K Rex et al. 'Colonoscopic miss rates of adenomas determined by back-to-back colonoscopies'. en. In: *Gastroenterology* 112.1 (Jan. 1997), pp. 24–28.
- [7] Alexander Robey, Hamed Hassani and George J. Pappas. *Model-Based Robust Deep Learning: Generalizing to Natural, Out-of-Distribution Data*. 2020. arXiv: 2005.10247 [cs.LG].
- [8] Dinggang Shen, Guorong Wu and Heung-Il Suk. 'Deep Learning in Medical Image Analysis'. In: *Annual Review of Biomedical Engineering* 19.1 (2017). PMID: 28301734, pp. 221–248. DOI: 10 . 1146 / annurev-bioeng-071516-044442. eprint: <https://doi.org/10.1146/annurev-bioeng-071516-044442>. URL: <https://doi.org/10.1146/annurev-bioeng-071516-044442>.
- [9] Sidney J. Winawer et al. 'Prevention of Colorectal Cancer by Colonoscopic Polypectomy'. In: *New England Journal of Medicine* 329.27 (1993). PMID: 8247072, pp. 1977–1981. DOI: 10 . 1056 / NEJM199312303292701. eprint: <https://doi.org/10.1056/NEJM199312303292701>.

NEJM199312303292701. URL: [https : / / doi . org / 10 . 1056 / NEJM199312303292701](https://doi.org/10.1056/NEJM199312303292701).

- [10] John R. Zech et al. 'Variable generalization performance of a deep learning model to detect pneumonia in chest radiographs: A cross-sectional study'. In: *PLOS Medicine* 15.11 (Nov. 2018), pp. 1–17. DOI: 10.1371/journal.pmed.1002683. URL: <https://doi.org/10.1371/journal.pmed.1002683>.