

Leveraging user-interaction and auxiliary data to learn from small data.

Romain Mormont

Systems and modeling unit,
Department of EE & CS,
University of Liège, Belgium

3rd October 2016

Context

- **Machine learning** (ML) has recently gained popularity through several successes (DeepMind, Google and Tesla self-driving cars, IBM Watson,...)
- Those were mostly possible thanks to:
 - **Plenty of data** (*big data*)
 - High computational power
- Those criterion are essentials to get the best performances out of common ML methods

Problem

- **In practice, those two criteria are not always met.**
- Especially, useful (i.e. labelled) data can be scarce. One can then talk about **small data** !

Small data: "the amount of data is not large enough with respect to the complexity of the task at hands"

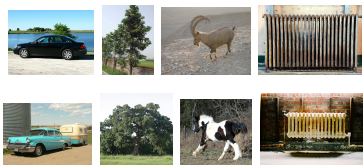
N.B.: in *small data* settings, the amount of data is not small in absolute terms

Small data: illustration

Big data ImageNet¹

Dataset:

- 14M labelled images
- Enough data to learn inter- and intra-class variabilities

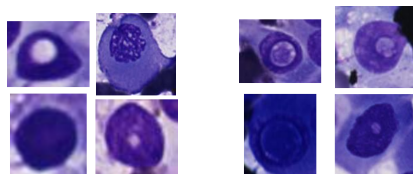


Typical images from ImageNet

Small data Thyroid nodule malignancy

Dataset:

- 6K labelled images
- Most important objects (i.e. malignant) are rare



(a) Healthy

(b) Malignant

¹Deng & al. *ImageNet: A Large-Scale Hierarchical Image Database*. In CVPR09, 2009.

Objectifs

Explore and develop new methods adapted to *small data* problems.

Tracks and research questions:

1. How to **integrate human operators to the learning process** in order to improve its performances ?
2. How to **leverage any available auxiliary data** for reaching the same goal ?

Methodology:

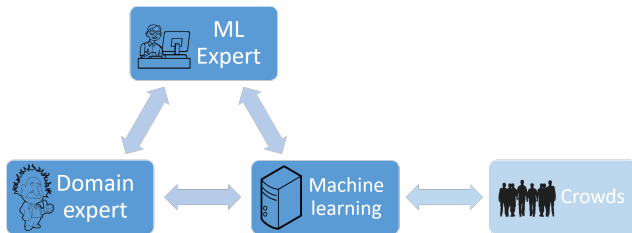
- Gather data
- Develop solution
- Test and improve it on benchmark problems
- Validate on real problems

Question 1: human in the loop ?

From ...



To ...



Question 1: human in the loop ?

State of the art:

- Mostly **active learning** (i.e. smart example labelling)
- Few developments around richer forms of feedback

Methodological challenges:

- How to **minimize the number of interaction** ? (the method/system must *ask the right questions*)
- How to **minimize time between each interaction** ? (the method/system must *be fast and reactive*)
- How to **integrate the feedbacks** to the learning process ?

Question 2: auxiliary data ?

TODO: Again with an illustration ? Transfert learning, imprecise/precise data, multi-modal images,...

Cas d'étude

Illustration cytomine

⇒ rare object detection and categorization in high-resolution tissue images

⇒ en production peut être utilisé pour collecter du feedback utilisateur

Merci pour votre attention !
Des questions ?

Backup slides

TODO:

- calendrier,
- résultats du TFE,
- applications à moyen ou long terme,
- publications récentes. . .
- théorie (ML,...)
- thyroid whole slide image for example illustration