CNN face à l'adversité autour d'un cas d'étude

J.E Campagne (IJCLab/Orsay) Action Dark Energy 2020









Outline

- Photometric redshift determination : the use case
- Case of *Inception* model first developed by J. Pasquet et al.
- Adversarial Samples
- Some results towards robustness
- Summary/Outlooks

Photo-z: photometric redshift



reference.

Methods for Photo-z

Since the pioneering work in the 60's, several methods have been developed to estimate the redshift from the multi-bands photometric measurements, basically:



• *Template* fitting

- eg. uses the SEDs of known galaxies and a fitting method
- Feature based Machine Learning
 - Uses a certain number of precooked features extracted from the measurements and feed to an engine as k-NN, NN/MLP, Decision Tree, BDT or Random Forest



Image based Deep Learning

Nb. Absolutely non exhaustive list of methods.

Possible combination

« Inception » for photo-z

Pasquet J., Bertin E., Treyer M., Arnouts S., Fouchez D., 2019, A&A, 621, A26 arXiv:1806.06607v2

Inspired from GoogleLeNet with multi-levels of conv-layers (Szegedy et al. 2014) 27.5M parameters Convolution (6.5% param.) FC (93.5% param.) ~30Layers Output Photo-z SDSS DR12 images + redenning 5x64x64 z=[0.02537716], ebv=[0.020828]

Results from my Inception



Zspec

Training/Test samples/plot :100k/100k/10k from a total of \sim 600k input dataset

Output: z-photo (regression, MSE)*

Bias ×10–4	<i>σmad</i> ×10−3	η(%)
+1.98	10.31	0.47

Results in agreement with J. Pasquet et al.

(*: arXiv:2002.10154 I was using the output of Pasquet et al.)

Common work

Teams spend some times to:

- Elaborate ML/DL architectures
- Apply some ML paradigm to tune the hyper-parameters using for instance: the triptych Training/Testing/Validation sets, Under/Over fitting aspects
- Compare their results against "State-of-the-art" competitors
- Perform systematics studies* on the Input Data: eg. are they representative of the use-case, what about their quality...

But, haven't we forgotten something?

(*: The study of J. Pasquet et al. is very detailed)

Adversarial samples: brief history

- After "**AlexNet**" the winner of ImageNet competition **2012**
- Topic rising since Szegedy et al. (2013): "Intriguing properties of neural networks"
- 1st explanation Goodfellow et al. (2014) : "Explaining and Harnessing Adversarial Examples"
- Part of the NIPS ' **2017** Competition
- Kurakin et al @ ICRL 2017: "Adversarial Machine Learning at Scale"
- Ilyas et al (2019): "Adversarial Examples Are Not Bugs, They Are Features"
- Madry et al (2017-19) @ ICLR 18: "Towards Deep Learning Models Resistant to Adversarial Attacks"
- ... Towards a deeper understanding of what is going on and how to overtake this intrinsic problem.

NSIP: Neural Information Processing Systems ICLR: International Conference on Learning Representations



Empirical risk/adversarial sample

 $\{x_i, z_i\}_{i \le N} \in D_{train}$ Eg. x_i : images, z_i : spectro-z



Min-max/saddle point problem: no general solution in non-convex problem
Which norm |δ|, which value of ε?

Simple perturbation mechanism

Goodfellow et al. 2014 Fast Sign Gradient Method f_{θ} "linear" + $\|\delta\|_{\infty} \leq \varepsilon \Rightarrow \delta^*(x) = \varepsilon \times \operatorname{sign}\left(\nabla_{\delta}\ell(f_{\theta}(x+\delta), z)\right)$



FSGM perturbations impact



(nb. <u>arXiv:2002.10154</u> I was using $\varepsilon = 10^{-2}$ which produce of course even more dramatic effects, but perturbation mode visible)

Some other results

- If one train 5 indepedant *Inception* models, and also a simpler CNN model, with each times a different set of training samples:
 - With *unperturbed* images: same results for each model (good), one can combine them.
 - With *perturbed images* build with one *Inception* model,
 - also perturbes the other *Inception* models: ie. **Combining different models cannot solve the problems**
 - Also perturbes different architecture models
- A perturbation δ^{*} is in principle directly linked to the original x unperturbed image, but it turns out that a single perturbation can impact a large number of images.
- One can use different ways to produce adversarial samples, I use here the FSGM method for the seak of clarity and also that one should first takle such simple perturbation.

What to do?

- A1: What's the problem ?
- A2: Bury one's head in the sand...
- A3: « These kind of perturbations will never append ! »: are you sure ? It is true that trying to generate the perturbations with « real artefact » is not so easy, at least faint objects can mimic small perturbations but not as efficient as the FSGM ones. The question is still open.
- Take it seriously as a sign of a certain (intrinsic) weakness:
 - Training ?
 - Architecture ?
 - Both ?

Countermeasures ?



Countermeasures ?

What about the training?

$$\theta_{t+1} = \theta_t - \alpha \frac{1}{|B_{train}|} \sum_{(x,z)\sim B_{train}} \nabla_{\theta} \left[\max_{\|\delta\| \leq \varepsilon} \ell(f_{\theta}(x+\delta), z) \right]_{\theta=\theta_t}$$

Solution not known in
the general case.
$$\nabla_{\theta} \ \ell(f_{\theta}(x+\delta^*), z) \right]$$
J. Danskin 1966
Convex case

$$\boldsymbol{\delta}^* = \max_{\|\boldsymbol{\delta}\| \leq \varepsilon} \ell(f_{\theta}(x + \boldsymbol{\delta}), z)$$

2) Mix up normal images & adversarial ones acts as regularisation terms.

1)

Finlay et al. 2018; Bietti et al. 2018

Adversarial training

During the training, add a certain fraction of adversarial samples. They act as a regularisation.

I love chemistry!

Training with 50% adv. images **FSGM** using $\epsilon=0.01$ (10x the attack)



Training	Images	Bias × 10–4	<i>σmad</i> ×10−3	η (%)
Train-Robust (Minℓtest ~4.4 10 ⁻⁴)	Non- perturbed	+0.67	14.27	1.95
	Perturbed	+5.54	15.56	2.14
Train-Classique (Min ℓ test ~2.2 10 ⁻⁴)	Non- perturbed	+1.98	10.31	0.47
	Perturbed	+66.40	33.63	5.31



Summary/outlooks

- The classical training/testing/validation triptych is not enough to guarantee the generalisation power of a network. Notice that the problem in more general than CNN (ie. DT, Gradient Boosted DT, R may also be affected as described in reference (Chen et al. 2019)).
- I've shown that mix up normal images with FSGM perturbed images gives some good results for *Inception* robustness, but for normal images the results are worse than those obtained with classical training.
- But this is **not the end of the story**: *Inception* is not immune against more aggressive methods. Some countermeasures have been elaborated but still it is a very active research domain as no satisfactory solution exists yet
- So, what next?
 - Change the architecture ? Well, I have tried several "classical" architecture but w/o real success, especially considering the robustness
 - Go back to the origin of the ML paradigm and modify the architecture with operators which do not need to be trained and are more stable against perturbations. Active research, new results will come, stay tuned!



Methods for Photo-z

Since the pioneering work in the 60's, several methods have been developed to estimate the redshift from the multi-bands photometric measurements, basically:



- template-fitting
 - Uses the SED and a method of fit
 - since Loh & Spillar 1986 ~30 galaxies in cluster 0024+1654,..., Beck et al 2016...
 - for LSST eg. Gorecki et al 2014 and Ansari et al 2019
- feature based Machine Learning
 - Uses a certain number of predefined features extracted from the measurements and feed to an engine as k-NN, NN/MLP, Decision Tree, BDT or Random Forest
 - Eg. Csabai et al. 2007 (k-NN) used by Beck et al 2016, Gorecki et al 2014 (NN) ,Ansari et al 2019 (BDT)...
- image based Deep Learning

Nb. Absolutely non exhaustive list of contributions.

Multi-steps perturbations

$$\begin{split} \delta^{*}(x) &\equiv \underset{\|\delta\| \leq \varepsilon}{\operatorname{argmax}} \ell(f_{\theta}(x+\delta), z) & \delta \leftarrow \delta + \underset{\|u\| \leq \alpha}{\operatorname{argmax}} \begin{bmatrix} u^{T} \cdot \nabla_{\delta} \ell(f_{\theta}(x+\delta), z) \end{bmatrix} \\ & \delta \leftarrow \delta + \underset{\|u\| \leq \alpha}{\operatorname{argmax}} \begin{bmatrix} u^{T} \cdot \nabla_{\delta} \ell(f_{\theta}(x+\delta), z) \end{bmatrix} \\ & \text{Kurakin et al 2016} \\ & \text{Kurakin et al 2016} \\ & \text{Charge terms of the sign} \left(\nabla_{\delta} \ell(f_{\theta}(x+\delta), z) \right) \end{bmatrix} \\ & \text{Charge terms of the sign} \left(\nabla_{\delta} \ell(f_{\theta}(x+\delta), z) \right) \end{bmatrix} \\ & \text{Charge terms of the sign} \left(\nabla_{\delta} \ell(f_{\theta}(x+\delta), z) \right) \end{bmatrix} \\ & \text{Charge terms of the sign} \left(\nabla_{\delta} \ell(f_{\theta}(x+\delta), z) \right) \end{bmatrix} \\ & \text{Charge terms of the sign} \left(\nabla_{\delta} \ell(f_{\theta}(x+\delta), z) \right) \end{bmatrix} \\ & \text{Charge terms of the sign} \left(\nabla_{\delta} \ell(f_{\theta}(x+\delta), z) \right) \end{bmatrix} \\ & \text{Charge terms of the sign} \left(\nabla_{\delta} \ell(f_{\theta}(x+\delta), z) \right) \end{bmatrix} \\ & \text{Charge terms of the sign} \left(\nabla_{\delta} \ell(f_{\theta}(x+\delta), z) \right) \end{bmatrix} \\ & \text{Charge terms of the sign} \left(\nabla_{\delta} \ell(f_{\theta}(x+\delta), z) \right) \end{bmatrix} \\ & \text{Charge terms of the sign} \left(\nabla_{\delta} \ell(f_{\theta}(x+\delta), z) \right) \end{bmatrix} \\ & \text{Charge terms of the sign} \left(\nabla_{\delta} \ell(f_{\theta}(x+\delta), z) \right) \end{bmatrix} \\ & \text{Charge terms of the sign} \left(\nabla_{\delta} \ell(f_{\theta}(x+\delta), z) \right) \end{bmatrix} \\ & \text{Charge terms of the sign} \left(\nabla_{\delta} \ell(f_{\theta}(x+\delta), z) \right) \end{bmatrix} \\ & \text{Charge terms of the sign} \left(\nabla_{\delta} \ell(f_{\theta}(x+\delta), z) \right) \end{bmatrix} \\ & \text{Charge terms of the sign} \left(\nabla_{\delta} \ell(f_{\theta}(x+\delta), z) \right) \end{bmatrix} \\ & \text{Charge terms of the sign} \left(\nabla_{\delta} \ell(f_{\theta}(x+\delta), z) \right) \end{bmatrix} \\ & \text{Charge terms of the sign} \left(\nabla_{\delta} \ell(f_{\theta}(x+\delta), z) \right) \end{bmatrix} \\ & \text{Charge terms of the sign} \left(\nabla_{\delta} \ell(f_{\theta}(x+\delta), z) \right) \end{bmatrix} \\ & \text{Charge terms of the sign} \left(\nabla_{\delta} \ell(f_{\theta}(x+\delta), z) \right) \end{bmatrix} \\ & \text{Charge terms of the sign} \left(\nabla_{\delta} \ell(f_{\theta}(x+\delta), z) \right) \\ & \text{Charge terms of the sign} \left(\nabla_{\delta} \ell(f_{\theta}(x+\delta), z) \right) \end{bmatrix} \\ & \text{Charge terms of the sign} \left(\nabla_{\delta} \ell(f_{\theta}(x+\delta), z) \right) \\ & \text{Charge terms of the sign} \left(\nabla_{\delta} \ell(f_{\theta}(x+\delta), z) \right) \\ & \text{Charge terms of the sign} \left(\nabla_{\delta} \ell(f_{\theta}(x+\delta), z) \right) \\ & \text{Charge terms of the sign} \left(\nabla_{\delta} \ell(f_{\theta}(x+\delta), z) \right) \\ & \text{Charge terms of the sign} \left(\nabla_{\delta} \ell(f_{\theta}(x+\delta), z) \right) \\ & \text{Charge terms of the sign} \left(\nabla_{\delta} \ell(f_{\theta}(x+\delta), z) \right) \\ & \text{Charge terms of the sign} \left(\nabla_{\delta} \ell(f_{\theta}(x+\delta), z) \right) \\ & \text{Charge terms of terms of terms of terms of terms$$

Some syst. studies

These are a very short summary of the J. Pasquet et al. thorough study.

Item	Comments
Galactic reddening (extra features added at the level of the FC part)	a strong reddening-dependent bias is observed If the information is not provided
Galaxy inclination	the CNN is very robust: large sample & data augmentation
Neighboring galaxies	The CNN learn how to improve redshift with neightboors at z>0.1
Variations throughout the surveyed area	Deviations in the SZ and Strip 82 of the SDSS dataset
PSF	induce a small but measurable amount of systematics on the estimated redshifts. Info can be added at the FC input (not done).

DL: what is promised?





Cascade of Convolutional layers

Pasquet J., Bertin E., Treyer M., Arnouts S., Fouchez D., 2019, A&A, 621, A26 <u>arXiv:1806.06607v2</u>

1) Variation: D'Isanto & Polsterer (2018) with a Gaussian Mixture Model as output

2) CNN architecture is used in other context: eg. g-g lens finding algo (Lanusse et al 2018), deblending (Burke et al 2019), objects classification (Gonzales et al 2018),...
2) 2

... Non exhaustive list !

Input images