Adversarial Training of Neural Networks against Systematic Uncertainty

Gilles Louppe New York University g.louppe@nyu.edu

Abstract

1 Introduction

[GL: Distinction between statistical and systematic uncertainty.] [GL: Define nuisance parameters.] [GL: We want to build an accurate classifier whose output remains invariant with respect to systematic uncertainties.] [GL: Motivate the criterion: allow to derive guarantees no matter λ , while f may otherwise be good or very bad depending on the nuisance if trained on the mixture. E.g. to form hypothesis to be later confirmed by data.]

2 Problem statement

Let assume a probability space (Ω, \mathcal{F}, P) , where Ω is a sample space, \mathcal{F} is a set of events and P is a probability measure. Let consider the multivariate random variables $X_{\lambda}: \Omega \mapsto \mathbb{R}^p$ and $Y: \Omega \mapsto \mathcal{Y}$, where X_{λ} depends on a nuisance parameter λ whose values define the family of its systematic uncertainties. That is, X_{λ} and Y induce together a joint probability distribution $p(X,Y|\lambda)$, where the conditional on λ denotes X_{λ} . For training, let further assume a finite set $\{x_i,y_i,\lambda_i\}_{i=1}^N$ of realizations $X_{\lambda_i}(\omega_i), Y(\omega_i)$, for $\omega_i \in \Omega$ and known values λ_i of the nuisance parameter. Our goal is to learn a function $f: \mathbb{R}^p \mapsto \mathcal{Y}$ (e.g., a classifier if \mathcal{Y} is a finite set of classes) minimizing the expected value of a loss $L(Y, f(X_{\lambda}))$, with the constraint that $f(X_{\lambda})$ should be robust to the value of the nuisance parameter λ – which remains unknown at test time. More specifically, we aim at building f such that in the ideal case

$$f(X_{\lambda_i}(\omega)) = f(X_{\lambda_i}(\omega)) \tag{1}$$

for any sample $\omega \in \Omega$ and any λ_i, λ_j pair of values of the nuisance parameter.

Since we do not have training tuples $(X_{\lambda_i}(\omega), X_{\lambda_j}(\omega))$ (for the same unknown ω), we propose instead to solve the closely related problem of finding a predictive function f such that

$$P(\{\omega|f(X_{\lambda_i}(\omega)) = y\}) = P(\{\omega'|f(X_{\lambda_i}(\omega')) = y\}) \text{ for all } y \in \mathcal{Y}.$$
 (2)

In words, we are looking for a predictive function f such that the distribution of $f(X_{\lambda})$ is invariant with respect to the nuisance parameter λ . Note that a function f for which Eqn. 1 is true necessarily satisfies Eqn. 2. The converse is however in general not true, since the sets of samples $\{\omega|f(X_{\lambda_i}(\omega))=y\}$ and $\{\omega'|f(X_{\lambda_j}(\omega'))=y\}$ do not need to be the same for the equality to hold.

3 Method

Adversarial training was first proposed by [1] as a way to build a generative model capable of producing samples from random noise $z \sim p_Z$. More specifically, the authors pit a generative model

29th Conference on Neural Information Processing Systems (NIPS 2016), Barcelona, Spain.

G against an adversary classifier D whose repelling objective is to recognize real from generated data. Both models G and D are trained simultaneously, in such a way that G learns to produce samples that are difficult to identify by D, while D incrementally adapts to changes in G. At the equilibrium, G models a distribution whose samples can be identified by D only by chance. In other words, assuming enough capacity in D and G, the distribution $p_{G(Z)}$ eventually converges towards the real distribution p_X .

In this work, we repurpose adversarial training for regularizing the construction of the predictive model f so as to satisfy Eqn. 2.

[GL: describe baseline] [GL: describe adversarial approach] [GL: proof that it solves Eqn. 2]

4 Experiments

5 Related work

[GL: Similar to domain adaptation, but with infinitely many domains, as parameterized by λ .]

6 Conclusions

Acknowledgments

References

[1] I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio, "Generative adversarial nets," in *Advances in Neural Information Processing Systems*, pp. 2672–2680, 2014.