

SYSTEMATIC AWARE LEARNING

a case study in High Energy Physics

Victor Estrade, Cecile Germain, Isabelle Guyon, David Rousseau

July 10, 2018

LRI, University of Paris-Sud, University of Paris-Saclay

Systematics

what are systematic effects and systematic uncertainties ?

Domain adaptation

learning to be robust against distribution shifts

Benchmark

enable fast experimentation

Experimental results

deep study of the proposed techniques

Conclusions and Perspectives

Systematics

what are systematic effects and systematic uncertainties ?

Domain apdaption

learning to be robust against distribution shifts

Benchmark

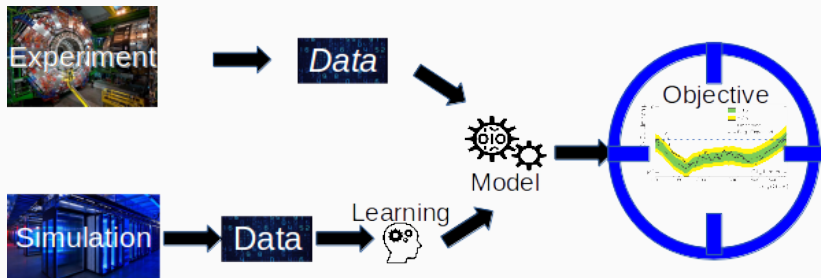
enable fast experimentation

Experimental results

deep study of the proposed technics

Conclusions and Perspectives

A COMMON SCENARIO IN EXPERIMENTAL SCIENCE



- Simulations are needed to understand data
- Machine learning is central in the pipeline

ERROR SOURCES

$$\text{measure} = \text{value} \pm \sigma_{\text{stat}} \pm \sigma_{\text{syst}}$$

Statistical error

- Empirical \neq asymptotic (Lack of data)
- Noise

Systematic error

"known unknowns"

- Apparatus imperfection
- Simulation imperfection
- Lack of theoretical knowledge
- Bugs

Model selection

- Limited model capacity
- Biased model choice
- Limited computation resources

EXAMPLE OF SYSTEMATIC EFFECT : CAMERA ROTATION

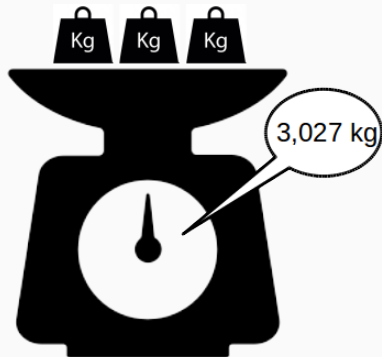
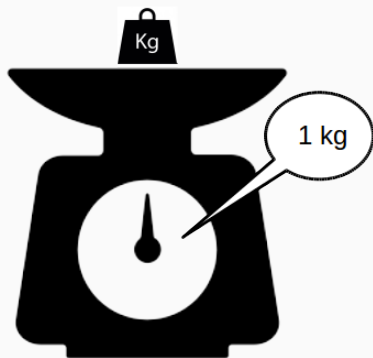
Original digits



Apparatus imperfection slightly rotates the image



EXAMPLE OF SYSTEMATIC EFFECT : SCALE PROBLEM



The scaling issue : true value = $(1 + z) \times \text{measured value}$

CHARACTERIZE SYSTEMATICS

Skewing function : $d(x, z)$

- Rotation of image input
- Rescaling

Nuisance parameter : z

- Angle of the rotation
- Scale factor

In real life there is several nuisance parameters with different impact on the data

Baby steps here : let's start with just one

Systematics

what are systematic effects and systematic uncertainties ?

Domain adaptation

learning to be robust against distribution shifts

Benchmark

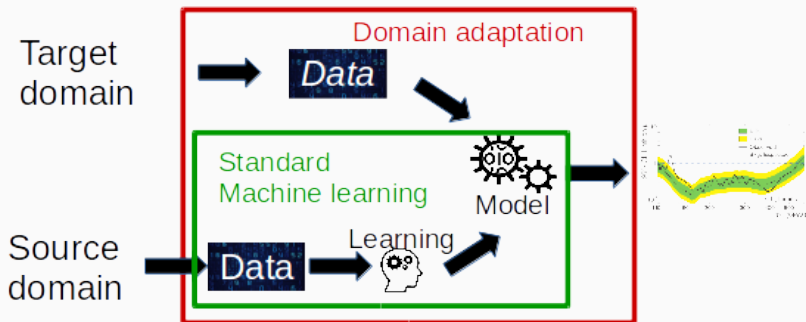
enable fast experimentation

Experimental results

deep study of the proposed techniques

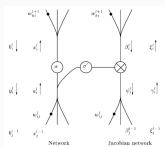
Conclusions and Perspectives

DOMAIN ADAPTATION (TRANSFER LEARNING)



Domain adaptation helps machine learning to be accurate on similar data

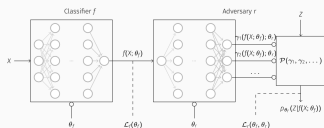
DEEP LEARNING ARCHITECTURES (A FEW SAMPLES)



Tangent Propagation (TP)

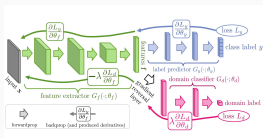
[Simard et al., 1991]

[Rifai et al., 2011]



Pivot Adversarial Network (PAN)

[Louppe et al., 2016]



Domain Adversarial Network (DAN)

[Ganin et al., 2015]

Generative adversarial networks (conceptual)



Generative Adversarial Network (GAN)

[Goodfellow et al., 2014]

TANGENT PROPAGATION IN A NUTSHELL

$$d(\text{3}, z) = \text{3}$$

- Regularize the derivative of the model according to the transformation.

$$loss = E_{standard} + \lambda \sum_{x \in Data} \left| \frac{\partial f(d(x, z); \theta)}{\partial z} \right|_{z=0}^2$$

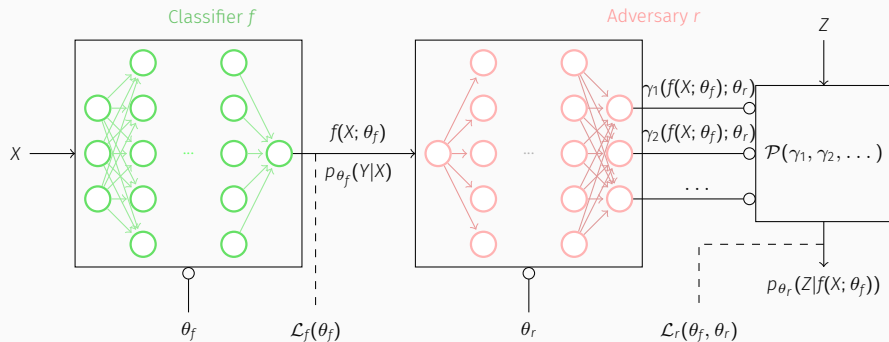
- Less data intensive than data augmentation
- Jacobian vector product trick : compute this derivative with a forward propagation through a "linearised" network.

$$\left. \frac{\partial f(d(x, z); \theta)}{\partial z} \right|_{z=0} = \nabla_x f(x; \theta) \cdot \left. \frac{\partial d(x, z)}{\partial z} \right|_{z=0}$$

[Simard et al., 1991] [Rifai et al., 2011]

PIVOT ADVERSARIAL NEURAL NETWORK

[Louppe et al., 2016]



- Learn the loss
- Makes it impossible to reconstruct Z from the output of the model
- Can include knowledge about nuisance parameter distribution

Systematics

what are systematic effects and systematic uncertainties ?

Domain apdaption

learning to be robust against distribution shifts

Benchmark

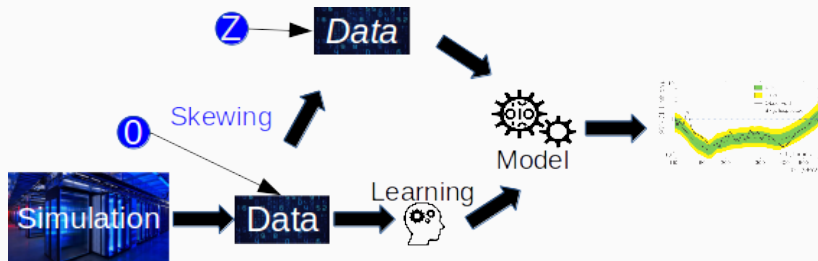
enable fast experimentation

Experimental results

deep study of the proposed technics

Conclusions and Perspectives

HEP BENCHMARK : SIMULATION DATA



- Simulation of the $H \rightarrow \tau\tau$ decay
- Nuisance parameter : τ energy scale ($\pm[1\%, 10\%]$)
- Data from HiggsML challenge [Adam-Bourdarios et al., 2014]¹
- Data from [Baldi et al., 2014]²

¹Available <http://opendata.cern.ch/record/328>

²Available <https://archive.ics.uci.edu/ml/datasets/HIGGSs>

ESTIMATE STATISTIC AND SYSTEMATIC ERROR

Final objective : measuring a cross section

Number of events N_z follow a Poisson distribution

The nuisance parameter induce a multiplicative error

After addition of the logarithmic derivatives we get :

$$\frac{\sigma_\mu}{\mu} = \sqrt{\left(\frac{\sqrt{s_0 + b_0}}{s_0}\right)^2 + \left(\frac{(s_z - s_0) + (b_z - b_0)}{s_0}\right)^2}$$

- $s = \sum_{S, score_i > t} w_i$, selected signals (True positives)
- $b = \sum_{B, score_i > t} w_i$, selected backgrounds (False positives)
- $*_0$, on the nominal data
- $*_z$, on the skewed data

Learning objective is to minimize the relative error $\frac{\sigma_\mu}{\mu}$

Systematics

what are systematic effects and systematic uncertainties ?

Domain adaptation

learning to be robust against distribution shifts

Benchmark

enable fast experimentation

Experimental results

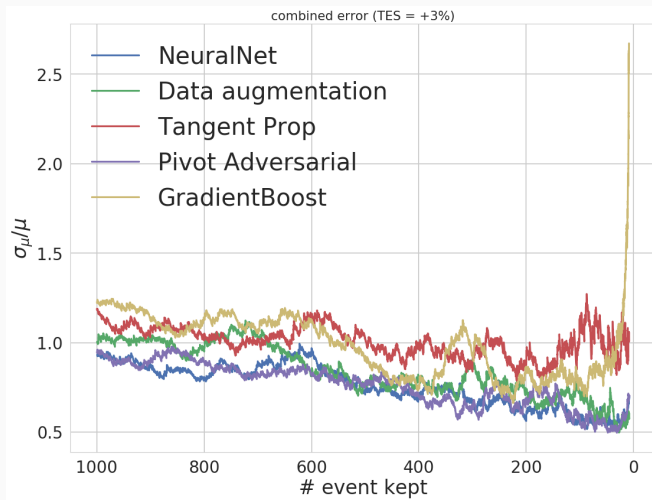
deep study of the proposed techniques

Conclusions and Perspectives

EXPERIMENTAL SETTINGS

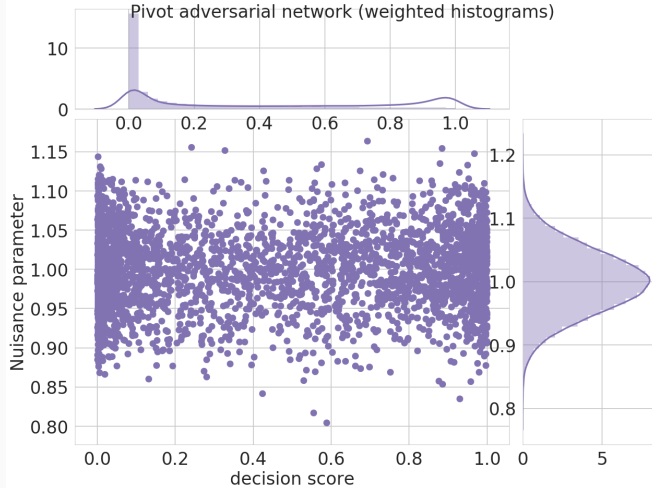
- Hyper parameters selected with grid search
- Models are compared with same structures (#neurons, non-linearities)
- Data augmentation and Pivot are fed with z drawn from a Gaussian distribution
- The others are trained with nominal data only
- Report estimated error along classification threshold

RESULTS OF THE CONTEST



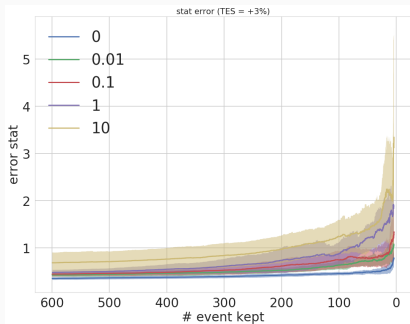
- Nothing beats the baseline (neural net)
- Tangent propagation is the worse.

MIXTURE FAILURE



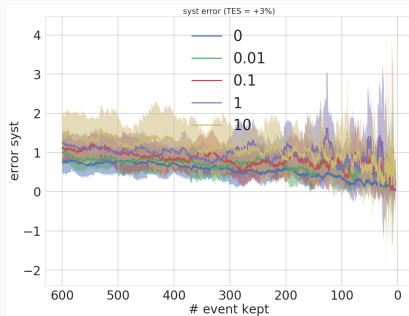
- Data augmentation & Pivot imitate a basic Neural Net
- Nothing to learn from the new skewed data instances ?

TANGENT PROPAGATION : LOOSING ON STAT ONLY



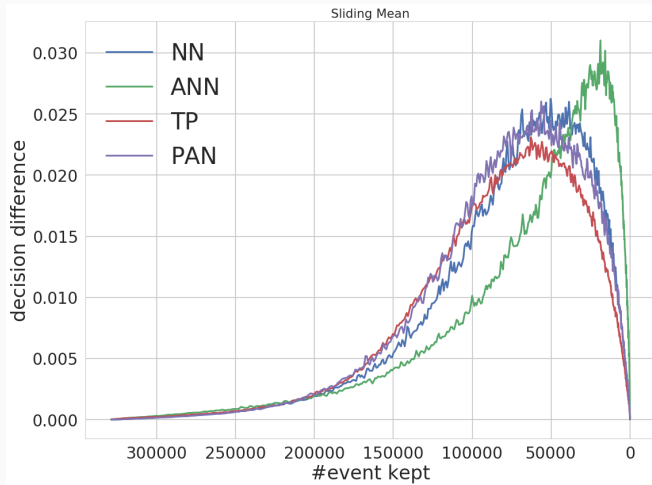
$$\frac{\sqrt{s_0 + b_0}}{s_0}$$

Performance loss mainly from statistical error



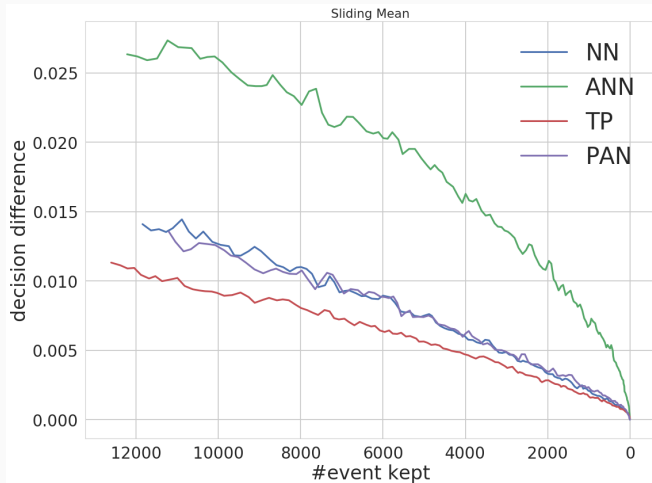
$$\frac{(s_z - s_0) + (b_z - b_0)}{s_0}$$

DECISION SHIFT VS SCORE (+3%)



- Tangent Propagation is reducing decision shift
- Mixture models are not imitating neural net
- Why $\frac{\sigma_{\mu}}{\mu}$ is not showing this behaviour ?

DECISION SHIFT VS SCORE (+3%)



- Tangent Propagation is reducing decision shift
- Mixture models are not imitating neural net
- Why $\frac{\sigma_{\mu}}{\mu}$ is not showing this behaviour ?

Systematics

what are systematic effects and systematic uncertainties ?

Domain adaptation

learning to be robust against distribution shifts

Benchmark

enable fast experimentation

Experimental results

deep study of the proposed techniques

Conclusions and Perspectives

Rank vs Score

- We've been tackling the problem in the wrong way.
- We don't need robustness along the classification score.
- We need the rank to be constant.

Perspective

- New toy dataset with controlled properties
- Explore domain adaptation giving robust ranking

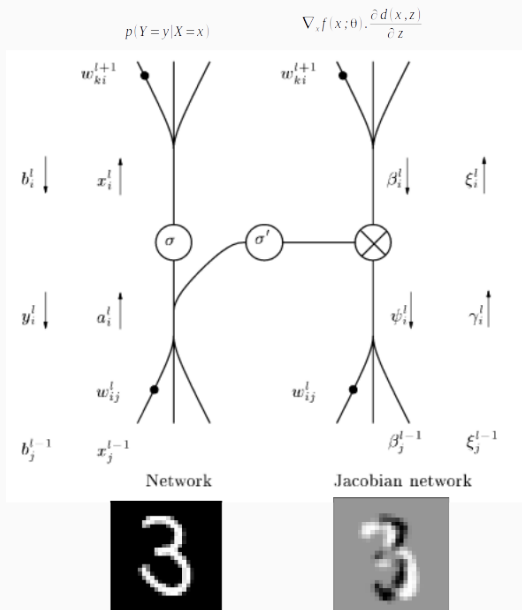
THANK YOU FOR YOUR ATTENTION

QUESTIONS ?

TANGENT VECTOR

$$\left. \frac{\partial d(\text{3}, z)}{\partial z} \right|_{z=0} = \frac{\text{3} - \text{3}}{2z} = \text{3}$$

JACOBIAN VECTOR PRODUCT TRICK



PIVOT ADVERSARIAL NEURAL NETWORK

Generative process $p(X, Y, Z)$

Train a neural net $f(X; \theta_f)$ to estimate the probability density $p(Y|X)$ (Z is marginalized)

We want to have the pivotal condition :

$$\forall(z, z'), p(f(X; \theta_f) = \text{score}|z) = p(f(X; \theta_f) = \text{score}|z')$$

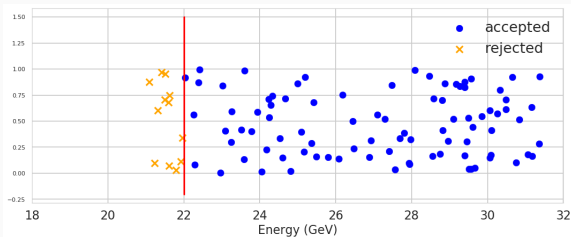
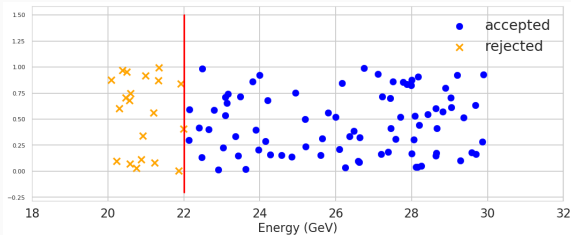
We want $f(X; \theta_f)$ and Z to be independent random variables

So we train an adversarial network to estimate $p(Z|f(X; \theta_f))$

[Louppe et al., 2016]

HEP BENCHMARK : TAU ENERGY SCALE

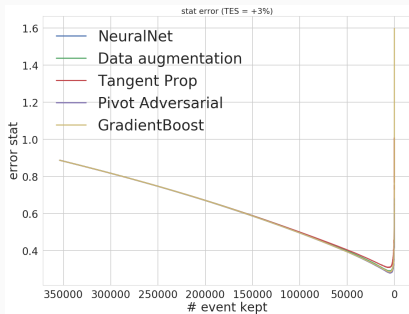
1. Scaling; 2. recompute derivate features; 3. cut data



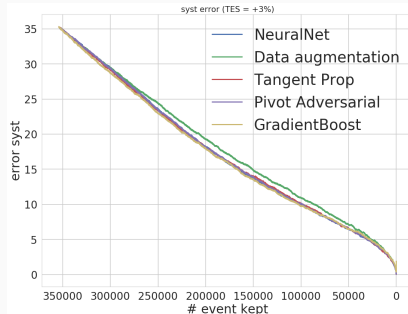
HEP BENCHMARK : MEASURING A CROSS SECTION

- Cross section = counting positive class
- $N_z = s_z + b_z$ selected event (Positives)
- $s_z = \sum_{S, score_i > t} w_i$, selected signals (True positives)
- $b_z = \sum_{B, score_i > t} w_i$, selected backgrounds (False positives)
- $\hat{s}_z = N_z - b_0$
- Normalized cross section $\mu_z = \frac{\hat{s}_z}{s_0} = \frac{s_z + b_z - b_0}{s_0}$

SYSTEMATICS DOMINATE



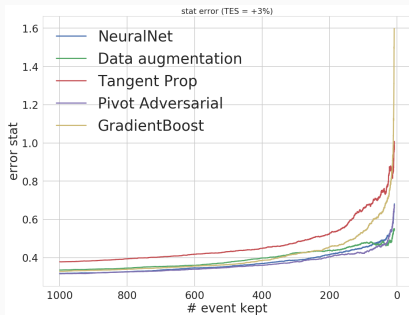
$$\frac{\sqrt{s_0 + b_0}}{s_0}$$



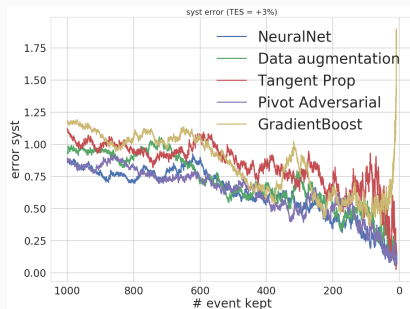
$$\frac{(s_z - s_0) + (b_z - b_0)}{s_0}$$

- The systematics dominates the measurement error

SYSTEMATICS DOMINATE



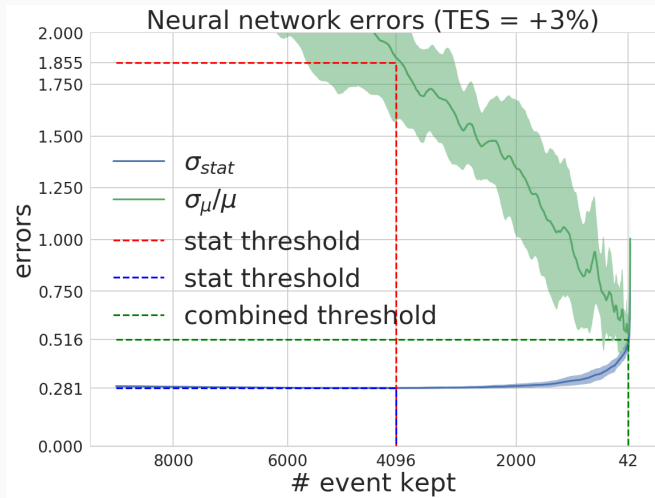
$$\frac{\sqrt{s_0 + b_0}}{s_0}$$



$$\frac{(s_z - s_0) + (b_z - b_0)}{s_0}$$

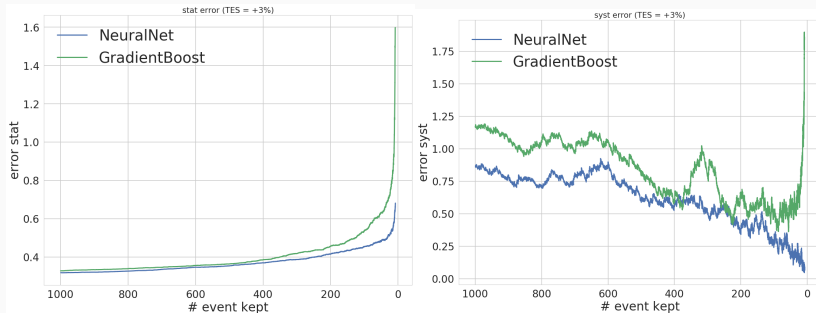
- The systematics dominates the measurement error
- But the statistical error cannot be ignored near minimum

THRESHOLD OPTIMIZATION



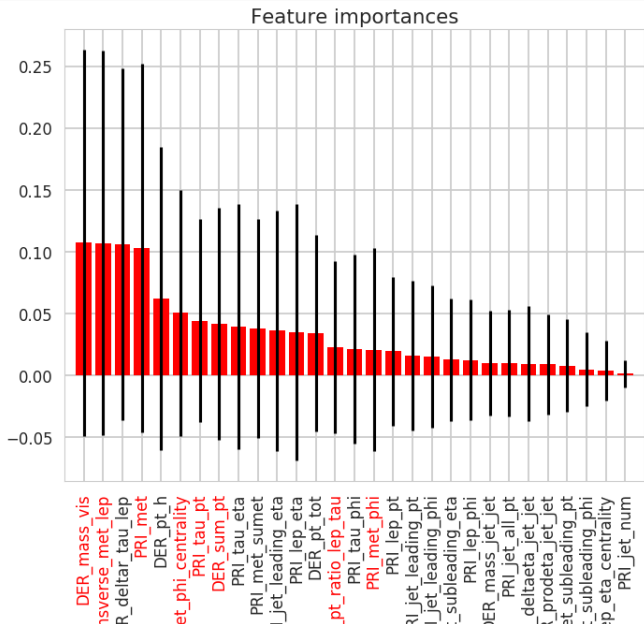
- Taking systematic effect into account drastically changes nb of event kept

GRADIENT BOOSTING : DETAILS

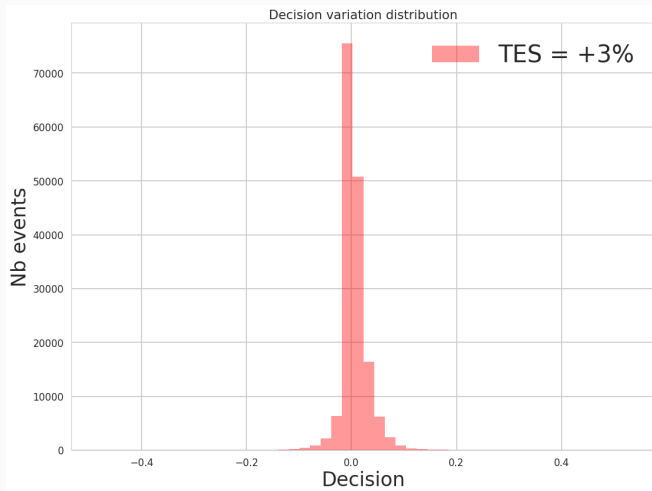


- sklearn 0.18
- 1000 trees
- maximum depth is 3 (i.e. 6 nodes maximum)
- Train only on nominal data ($z = 0$)

GRADIENT BOOSTING : USING SKEWED FEATURES ?

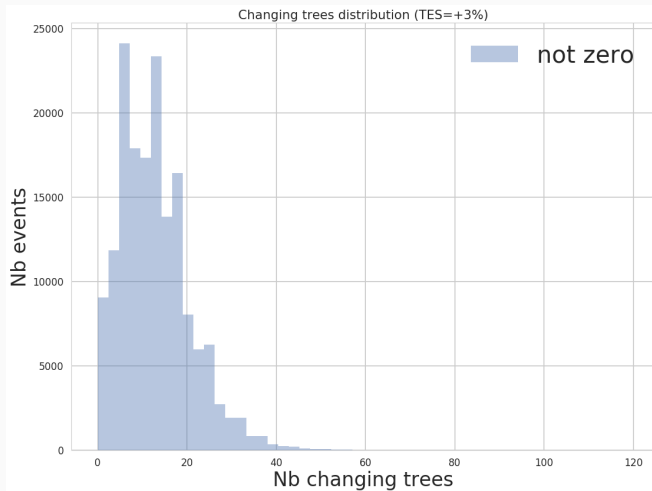


GRADIENT BOOSTING : SMALL DECISION VARIATION



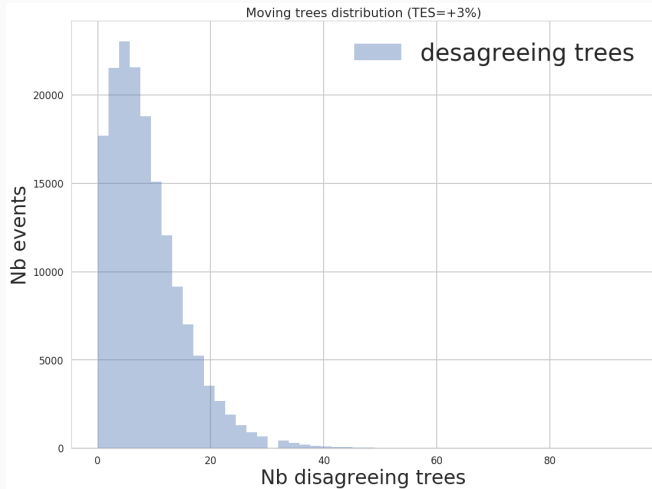
Score variation are small and goes both way

GRADIENT BOOSTING : FEW TREES ARE AFFECTED



Only a few trees among the 1000 are changing

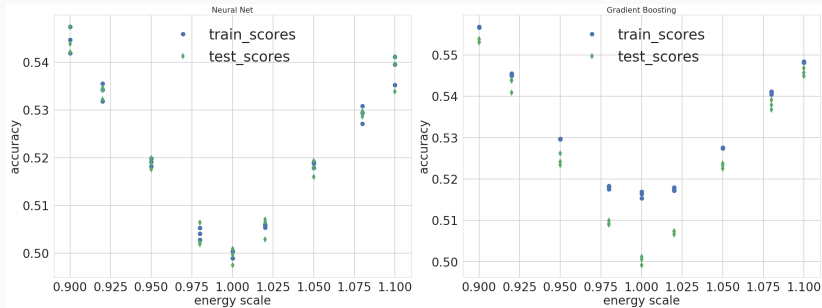
GRADIENT BOOSTING : TREE DISAGREEMENT



GRADIENT BOOSTING : CONCLUSION

- Based on small trees : decision function is constant almost everywhere
- Many trees disagree among themselves

HIGGS : SEPARABILITY ISSUE

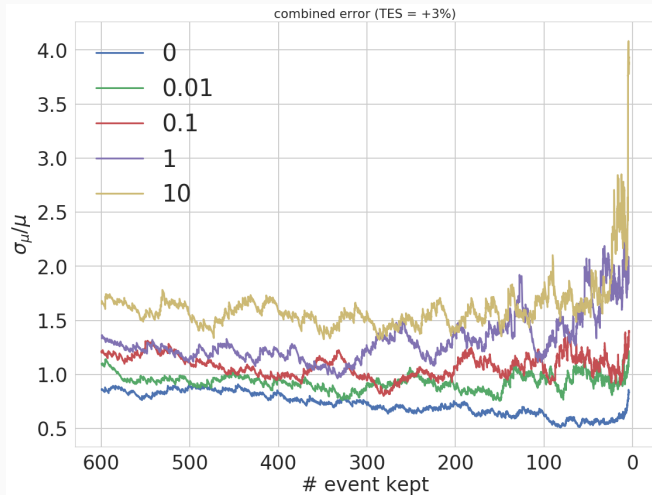


- low \mathcal{H} -divergence [Ben-David et al., 2010]

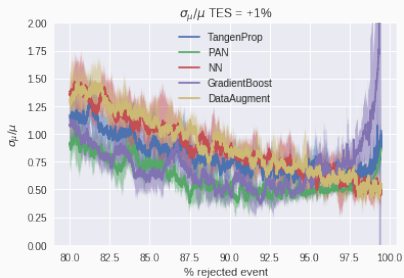
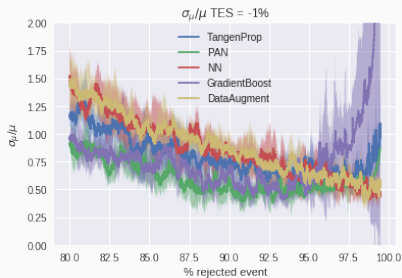
TANGENT PROPAGATION : DETAILS

- 3 hidden layers
- 120 neurons each
- Adam optimizer
- Train on nominal data + tangent vectors

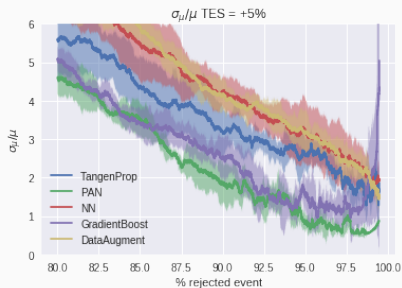
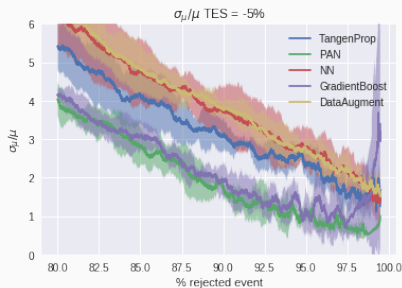
TANGENT PROPAGATION : LOOSING FOR EVERY λ



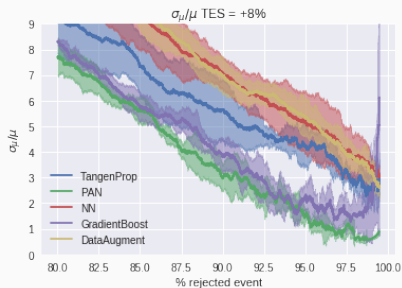
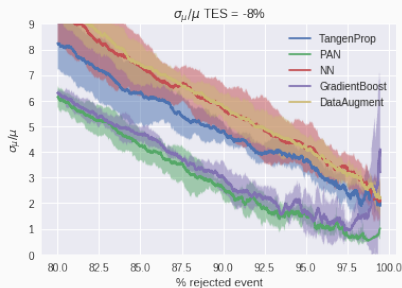
SYMMETRY ? (1%)



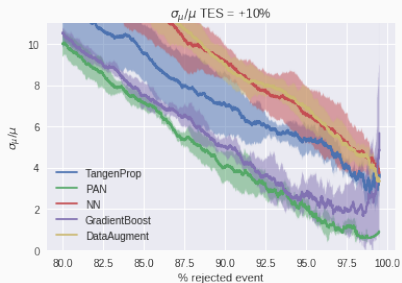
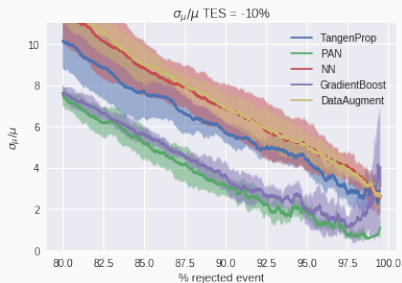
SYMMETRY ? (5%)



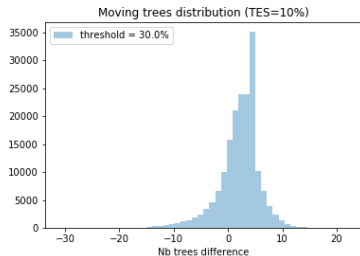
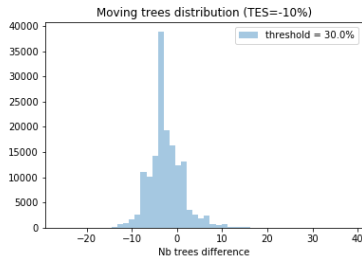
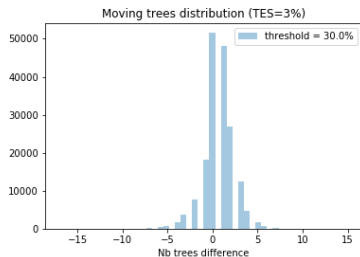
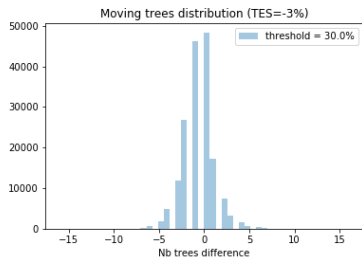
SYMMETRY ? (8%)



SYMMETRY ? (10%)

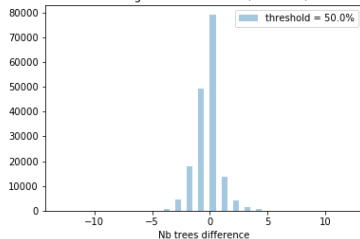


GRADIENT BOOSTING : DISAGREEMENT SYMMETRY

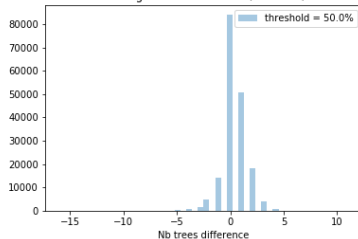


GRADIENT BOOSTING : STRONG DISAGREEMENT SYMMETRY

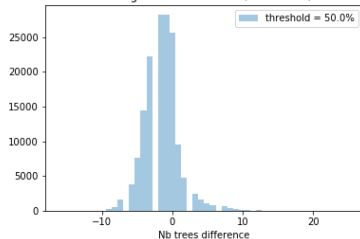
Moving trees distribution (TES=-3%)



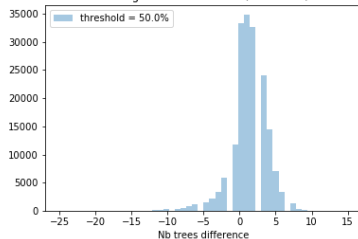
Moving trees distribution (TES=3%)



Moving trees distribution (TES=-10%)



Moving trees distribution (TES=10%)



invariant vs profiling.

Being invariant is too hard. In the end the nuisance params take only one value.

Ask David : "In practice, a NP is η and p_T dependent and affect each events differently" So in the end the tau energy scale may differ a little bit event wise ? Answer : No, the NP is the same for all the events but its impact on each event depend on other variables

The DER features may contribute to the robustness of GB ?



Adam-Bourdarios, C., Cowan, G., Germain, C., Guyon, I., Kégl, B., and Rousseau, D. (2014).

The Higgs boson machine learning challenge.

In *NIPS 2014 Workshop on High-energy Physics and Machine Learning*, volume 42 of *JMLR: Workshop and Conference Proceedings*, pages 19–55, Montreal, Canada.



Baldi, P., Sadowski, P., and Whiteson, D. (2014).

Enhanced Higgs to $\tau^+\tau^-$ Searches with Deep Learning.



Ben-David, S., Blitzer, J., Crammer, K., Kulesza, A., Pereira, F., and Vaughan, J. W. (2010).

A theory of learning from different domains.

Machine Learning, 79(1-2):151–175.



Ganin, Y., Ustinova, E., Ajakan, H., Germain, P., Larochelle, H., Laviolette, F., Marchand, M., and Lempitsky, V. (2015).

Domain-Adversarial Training of Neural Networks.

arXiv:1505.07818 [cs, stat].

arXiv: 1505.07818.



Goodfellow, I. J., Pouget-Abadie, J., Mirza, M., Xu, B., Warde-Farley, D., Ozair, S., Courville, A., and Bengio, Y. (2014).

Generative Adversarial Networks.

arXiv:1406.2661 [cs, stat].

arXiv: 1406.2661.



Louppe, G., Kagan, M., and Cranmer, K. (2016).

Learning to Pivot with Adversarial Networks.

arXiv:1611.01046 [physics, stat].

arXiv: 1611.01046.



Rifai, S., Dauphin, Y., Vincent, P., Bengio, Y., and Muller, X. (2011).

The Manifold Tangent Classifier.

In *NIPS*, volume 271, page 523.



Simard, P. Y., Victorri, B., LeCun, Y., and Denker, J. S. (1991).

Tangent Prop - A Formalism for Specifying Selected Invariances in an Adaptive Network.

In Moody, J. E., Hanson, S. J., and Lippmann, R., editors, *NIPS*, pages 895–903. Morgan Kaufmann.