SISSA MEETS ENS

2018-09-28

MACHINE LEARNING FOR HIGH ENERGY PHYSICS

Andrea De Simone





andrea.desimone@sissa.it

> The APP group at SISSA

interdisciplinary group working at interface of **particle physics**, **astrophysics and cosmology**

Address **fundamental issues** about our Universe: origin and evolution, nature of gravity, properties of dark matter and dark energy.

Pls in the group: A. De Simone, S. Liberati, P. Ullio, M. Viel

Pls affiliated from other SISSA groups: C. Baccigalupi, R. Percacci, S. Petcov, A. Romanino, P. Salucci

Pls affiliated form other institutions P. Creminelli, E. Sefusatti

Speech recognition





Recommender systems



Creative Paintings

Which of these images were created by a machine?



[Elgammal et al -1706.07068]

Generating Faces

These people do not exist!



A. De Simone

Autonomous driving

43 мрн



Games



Google's DeepMind plays Breakout



AlphaGo beats world champion Go 4-1

2010-05-06, 2:45pm



2010-05-06, 2:45pm Flash Crash!



lost ~9% in 36 minutes!



1. Machine Learning in Science

2. Open problems in High-Energy Physics (HEP)

3. Statistical test of dataset compatibility

4. Applications to HEP



1. Machine Learning in Science

2. Open problems in High-Energy Physics (HEP)

3. Statistical test of dataset compatibility

4. Applications to HEP





"A computer program is said to *learn* from **experience E** with respect to some class of **tasks T** and **performance measure P**,

if its performance at tasks in T, as measured by P, improves with experience E."

[Mitchell - 1997]

Why is Machine Learning so cool?

- a guide through big data (data mining)
- many diverse applications (from engineering to commerce to science)
- can help making our life better/easier
- ...
- can help scientists to do science better and faster



- More powerful machines

both speed and storage

- More data

almost everything is recorded!

- Easier access

internet revolution, easier to share big data

Supervised Learning



Unsupervised Learning

Reinforcement Learning

> Supervised Learning

. . .



Logistic Regression Neural Networks Decision Trees Nearest Neighbors Polynomial Regression Neural Networks Support Vector Machines Nearest Neighbors

machine "learns" the model

f(x) = y

> Unsupervised Learning

. . .





Cluster Analysis Dimensionality Reduction Anomaly Detection

machine "learns" patterns, structures, representations, etc. of the data

> ML in Medical Science

deep learning achieved state-of-the art results



[Litjens et al. - Medical Image Analysis 2017]

> ML in Medical Science

Dermatologist-level classification of skin cancer with deep neural networks



nature

LESIONS LEARN

> ML in Chemistry

Neural Networks for the prediction of organic chemistry reactions



predict probability of 17 different reaction types

- 0. Null Reaction
- 1. Nucleophilic substitution
- 2. Elimination
- 3. Nucleophilic Substitution with Methyl Shift
- 4. Elimination with methyl shift
- 5. Hydrohalogenation (Markovnikov)
- 6. Hydrohalogenation (Anti-Markovnikov)
- 7. Hydration (Markovnikov)
- 8. Hydration (Anti-Markovnikov)
- 9. Alkoxymercuration-demercuration
- 10. Hydrogenation
- 11. Halogenation
- 12. Halohydrin formation
- 13. Epoxidation
- 14. Hydroxylation
- 15. Ozonolysis
- 16. Polymerization



(~85% accuracy on test set)

[Wei, Duvenaud, Aspuru-Guzik, 2016]

Detect extreme weather using deep learning

classification accuracy

Event Type	Train	Test	Train
			time
Tropical Cyclone	99%	99%	$\approx 30 \min$
Atmospheric River	90.5%	90%	6-7 hour
Weather Front	88.7%	89.4%	$\approx 30 \min$













> ML in Physics

Machine learning phases of matter



triangular-lattice Ising model: Tc/J = 3.65 ± 0.01 (exact: 3.64095...)



SICS

> ML in Physics



A. De Simone



1. Machine Learning in Science

2. Open problems in High-Energy Physics (HEP)

3. Statistical test of dataset compatibility

4. Applications to HEP

> The Standard Model



Standard Model of Elementary Particles

> Higgs Discovery

2012-07-04





Is there anything left to discover?

> Need to go beyond...



> Large Hadron Collider (LHC)

The most complex (and expensive) experiment ever built!

Cost ~ 4 GigaEur



counter-rotating proton beams in 27km circumference ring



detectors at 4 collision points:

- CMS - LHCb - ALICE

- ATLAS

> ATLAS experiment

~ 2100 physicists 37 countries 167 universities/labs







A. De Simone

> Signal vs Background



> Standard Analysis Pipeline



> Standard Analysis Pipeline



> New Physics ?

Searches for New Physics Beyond the Standard Model have been negative so far...



1. New Physics is not accessible by LHC

new particles are too light/heavy or interacting too weakly

2. We have not explored all the possibilities

new physics may be buried under large bkg or hiding behind unusual signatures

"Don't want to miss a thing" (in data)

closer look at currently available data get ready for upcoming data from next Run of LHC

Model-independent search

searches for specific models may be insensitive to unexpected / unknown / anomalous processes



1. Machine Learning in Science

2. Open problems in High-Energy Physics (HEP)

3. Statistical test of dataset compatibility

4. Applications to HEP

> New Statistical Test

[DS, Jacques - 1807.06038]

Want a statistical test for NP which is:

1. model-independent:

no assumption about underlying physical model to intepret data

more general

2. non-parametric:

compare two samples as a whole (not just their means, etc.)

fewer assumptions, no max likelihood estim.

3. un-binned:

high-dim feature space partitioned without rectangular bins retain full multi-dim info of data

[a.k.a. "homogeneity test"]

Two sets:
$$\mathcal{T} = \{x_1, \dots, x_{N_T}\} \stackrel{\text{iid}}{\sim} p_T$$
Benchmark: $\mathcal{B} = \{x'_1, \dots, x'_{N_B}\} \stackrel{\text{iid}}{\sim} p_B$

$$oldsymbol{x}_i,oldsymbol{x}_i'\in\mathbb{R}^L$$

probability distributions p_{B} , p_{T} unknown





Two sets:

Trial:
$$\mathcal{T} = \{x_1, \dots, x_{N_T}\} \stackrel{\text{iid}}{\sim} p_T$$
 $x_i, x'_i \in \mathbb{R}^D$ Benchmark: $\mathcal{B} = \{x'_1, \dots, x'_{N_B}\} \stackrel{\text{iid}}{\sim} p_B$ $x_i, x'_i \in \mathbb{R}^D$

probability distributions p_B, p_T unknown

Are *B*,*T* drawn from the same prob. distribution?



Benchmark and Trial samples

Two sets:

Trial:
$$\mathcal{T} = \{x_1, \dots, x_{N_T}\} \stackrel{\text{iid}}{\sim} p_T$$
 $x_i, x'_i \in \mathbb{R}^D$ Benchmark: $\mathcal{B} = \{x'_1, \dots, x'_{N_B}\} \stackrel{\text{iid}}{\sim} p_B$ $x_i, x'_i \in \mathbb{R}^D$

probability distributions p_{B} , p_{T} unknown

Are *B*,*T* drawn from the same prob. distribution?



A. De Simone

Why is it important?

- . . .

- decide whether two datasets can be analyzed jointly
- find anomalous data points (outliers)
- detect changes in the underlying distributions
- detect events in streams of data (time-series data)
- check if data are compatible with expectations



1. Density Estimator

reconstruct PDFs from samples

2. Test Statistic (TS)

measure "distance" between PDFs

3. TS distribution

→ associate probabilities to TS under null hypothesis H_0 : $p_B = p_T$

4. *p* -value

 \longrightarrow accept/reject H_0



Divide the space in squared bins?

- easy
 can use simple statistics (e.g. χ^2)
- hard/slow/impossible in high-D

Need un-binned multivariate approach

Find PDFs *estimators*: $\hat{p}_B(\boldsymbol{x}), \hat{p}_T(\boldsymbol{x})$ e.g. based on densities of points:

$$\hat{p}_{B,T}(\boldsymbol{x}) = \frac{\rho_{B,T}(\boldsymbol{x})}{N_{B,T}}$$

Nearest Neighbors!

```
[Schilling - 1986][Henze - 1988]
[Wang et al. - 2005,2006]
[Dasu et al. - 2006][Perez-Cruz - 2008]
[Sugiyama et al. - 2011][Kremer et al, 2015]
```



- Fix integer K.
- Choose query point *x_j* in *T* and draw it in *B*.



- Fix integer K.
- Choose query point *x_j* in *T* and draw it in *B*.
- Find the distance $r_{j,B}$ of the Kth-NN of x_j in *B*.







- Fix integer K.
- Choose query point *x_j* in *T* and draw it in *B*.
- Find the distance $r_{j,B}$ of the Kth-NN of x_j in *B*.
- Find the distance $r_{j,T}$ of the Kth-NN of x_j in *T*.





- Fix integer K.
- Choose query point *x_j* in *T* and draw it in *B*.
- Find the distance *r_{j,B}* of the Kth-NN of *x_j* in *B*.
- Find the distance $r_{j,T}$ of the Kth-NN of x_j in *T*.
- Estimate PDFs:

 $\hat{p}_T(\boldsymbol{x}_j)$

$$\hat{p}_B(oldsymbol{x}_j)$$
 :

$$= \frac{K}{N_B} \frac{1}{\omega_D r_{j,B}^D}$$
$$= \frac{K}{N_T - 1} \frac{1}{\omega_D r_j^2}$$

> 2. Test Statistic

- Measure of the "distance" between 2 PDFs
- Define Test Statistic: (detect under-/over-densities)

$$TS(\mathcal{T}) \equiv \frac{1}{N_T} \sum_{j=1}^{N_T} \log \frac{\hat{p}_T(\boldsymbol{x}_j)}{\hat{p}_B(\boldsymbol{x}_j)}$$

- Related to Kullback-Leibler divergence as: $\operatorname{TS}(\mathcal{T}) = \hat{D}_{\operatorname{KL}}(\hat{p}_T || \hat{p}_B)$ $D_{\operatorname{KL}}(p || q) \equiv \int_{\mathbb{R}^D} p(\boldsymbol{x}) \log \frac{p(\boldsymbol{x})}{q(\boldsymbol{x})} d\boldsymbol{x}$
- From NN-estimated PDFs: $TS(\mathcal{T}) = \frac{D}{N_T} \sum_{j=1}^{N_T} \log \frac{r_{j,B}}{r_{j,T}} + \log \frac{N_B}{N_T 1}$
- **Theorem:** this estimator converges to $D_{KL}(p_B || p_T)$, in large sample limit [Wang et al. - 2005,2006]

> 3. Test Statistic Distribution

How is TS distributed? **Permutation test!**

Assume $p_B = p_T$. Union set: $\mathcal{U} = \mathcal{T} \cup \mathcal{B}$



Distribution of TS under H_0 : $f(TS|H_0) \leftarrow {TS_n}^{\xi}$ [asymptotically normal with zero mean]

TS

> 4. p-value

• $\hat{\mu}, \hat{\sigma}$: mean, variance of TS distribution $f(TS|H_0)$

• Standardize the TS:
$$TS \rightarrow TS' \equiv \frac{TS - \hat{\mu}}{\hat{\sigma}}$$

• TS' distributed according to $f'(TS'|H_0) = \hat{\sigma}f(\hat{\mu} + \hat{\sigma}TS'|H_0)$



A. De Simone

> Gaussian Example



> Where are the discrepancies?

Bonus: Characterize regions with significant discrepancies

1. "Score" field over T: $Z(x_j) \equiv \frac{u(x_j) - \bar{u}}{\sigma_u}$



2. Identify points where Z(x) > cThey contribute the most to large TS_{obs} \rightarrow high-discrepancy (anomalous) regions

3. Apply a clustering algorithm to group them

with:
$$u(\boldsymbol{x}_j) \equiv \log \frac{r_{j,B}}{r_{j,T}}$$

TS_{obs} = $D \bar{u} + \text{const}$



> NN2ST: Summary

INPUT:

Trial sample: $\mathcal{T} = \{x_1, \dots, x_{N_T}\} \stackrel{\text{iid}}{\sim} p_T$ Benchmark sample: $\mathcal{B} = \{x'_1, \dots, x'_{N_B}\} \stackrel{\text{iid}}{\sim} p_B$ K:number of nearest neighborsNperm:number of permutations

 $oldsymbol{x}_i,oldsymbol{x}_i'\in\mathbb{R}^D$ $oldsymbol{p}_{\mathcal{B}},oldsymbol{p}_{\mathcal{T}}$ unknown

OUTPUT:

p-value of the null hypothesis H_0 : $p_B = p_T$

[check compatibility between 2 samples]

> NN2ST: Summary



Python code: <u>github.com/de-simone/NN2ST</u> Paper: [DS, Jacques - 1807.06038]

> NN2ST: Summary

- ✓ general, model-independent
- \checkmark fast, no optimization
- sensitive to unspecified signals
- ✓ useful when no variable can separate sig/bkg
- ✓ helps finding signal regions, optimal cuts, ...
- need to run for each sample pair
- permutation test is bottleneck
- Imited by sample accuracies



1. Machine Learning in Science

2. Open problems in High-Energy Physics (HEP)

3. Statistical test of dataset compatibility

4. Applications to HEP

> Standard Analysis Pipeline



> Our Method



> DM search @ LHC



- "proof-of-principle" study
- bkg: $Z \to \nu \bar{\nu} + (1,2)\, j$ sub-leading bkgs not included
- no full detector effects



Benchmark: BKG₁ Trials: BKG₂ + SIG K = 5 $N_{perm} = 3000$ 8 features:

- n. of jets
- p_T, η of 2 leading jets
- - E_T^{miss}, H_T

$$\Delta \phi_{E_T^{\mathrm{miss}},j_1}$$

> DM search @ LHC

BKG₁ (20k events)
T1: BKG₂ (20k events) + SIG₁ (2010 events)
T2: BKG₂ (20k events) + SIG₂ (375 events)
T3: BKG₂ (20k events) + SIG₃ (59 events)

$$N_{\rm sig} = N_B \times \frac{\sigma_{\rm signal}}{\sigma_{\rm bkg}}$$

Sample	Mz'	σ signal	Z	 in real world: expect degradation of results 	
T1	1.2 TeV	20.4 pb	>15 <i>o</i>	(uncertainties)	
T2	2 TeV	3.8 pb	10 σ	 the method has value, it is worth exploring 	
T3	3 TeV	0.6 pb	0.13 <i>σ</i>		

Expt. Collab. at CERN interested in applying this test

> DM search @ LHC





Directions for future work:

- inclusion and impact of uncertainties
- adaptive choice of K
- identifying discrepant regions in realistic situations (with *Z*-score method)
- validation tool for bkg: compatibility between MC sims. and data in control regions
- scalability

1. Golden age of Machine Learning

Big Data are everywhere

2. Innovative Statistical Test for New Physics

- Powerful and model-independent discovery tool
- Guidance for experimental searches

3. ML for science in germinal stage

Pioneering developments waiting ahead!