Machine-Learned Exclusion Limits without Binning

Rosa María Sandá Seoane

Instituto de Física Teórica UAM-CSIC

SATURNALIA 2022 December 21, 2022

э

Introduction

- About ML...
- Motivation
- ML Likelihood method
- Kernel Density Estimators

2 Applications

- Toy Model: Multivariate Gaussian distributions
- Search for a heavy SSM Z' in dilepton final states at the HL-LHC

3 Conclusions

About ML... Motivation ML Likelihood method Kernel Density Estimators

What is Machine Learning?

Machine Learning (ML) is the study of computer algorithms capable of building a mathematical model out of a data sample, by learning from examples.



The algorithm builds a predictive model without being explicitly programmed to do so

э

About ML... Motivation ML Likelihood method Kernel Density Estimators





None of the systems we have nowadays are real AI! The brain learns so efficiently that no ML method can match it.

э

About ML... Motivation ML Likelihood method Kernel Density Estimators

Learning Paradigms



About ML... Motivation ML Likelihood method Kernel Density Estimators

Supervised Learning

Given some labeled data $D = \{(\vec{x_1}, \vec{t_1}), ..., (\vec{x_n}, \vec{t_n})\}$ with features $\{\vec{x_i}\}$ and targets $\{\vec{t_i}\}$, the algorithm finds a mapping $\vec{t_i} = F(\vec{x_i})$



프 🕨 🖉

About ML... Motivation ML Likelihood method Kernel Density Estimators

Supervised Learning

Classification: $\{\vec{t_1}, ..., \vec{t_n}\}$ (finite set of labels) Regression: $\vec{t_i} \in \mathbb{R}^n$



(Linear Regression is the oldest ML algorithm!)

э

About ML... Motivation ML Likelihood method Kernel Density Estimators

Supervised Learning in HEP: examples

Classification: Jet Tagging (infering number of quarks and gluons inside it, or *prongs*)



 $\{\vec{x_i}\}$: jet mass, jet's tranverse momentum, 17 N-subjetiness variables. $\{\vec{t_i}\}=\{0,1\}$

(J.A. Aguilar-Saavedra, E. Arganda, F.R. Joaquim, J. Seabra, RMSS)

프 🕨 🖉

About ML... Motivation ML Likelihood method Kernel Density Estimators

Method: Machine-Learned Likelihood

E. Arganda, X. Marcano, V. Martín Lozano, A. D. Medina, A. D. Perez, M. Szewc, A. Szynkman

Eur. Phys. J. C **82**, no.11, 993 (2022)

E. Arganda, M. de los Rios, A. D. Perez, RMSS *PoS* ICHEP2022 (2022) 1226 E. Arganda, M. de los Rios, A. D. Perez, RMSS arXiv: 2211.04806

Machine-Learned Exclusion Limits without Binning

Ernesto Arganda ^{a,b} Andres D. Perez, ^b Martin de los Ries^{a,c} and Rosa María Sandá Scoane ^c]

¹⁰ Institute de Fisica Tetolos UAM-CSIC, C/ Nicolas Calerna 19-15, Compas de Cantolibano, 2010, Madrid, Spain ² BTLP, CONSCET - Dyte, de Fisica, Universidad Nacional de La Plata, C.C. 67, 1990 La Plata, Argentina ¹⁰ Deventuement de Fisica Tetrica, Volomidad Autónema de Madrid.

Espartanenzo de Faiza Zeorea, Consensada Atronana de Madro E-20049 Cantoblanco, Madrid, Spaix

.

A method for approximating optimal statistical significances with machine-learned likelihoods

Ernesto Arganda,^{e.b.}, Xabier Marcano,^{e.c.} Víctor Martín Lozano,^{d.e.} Anibal D. Medina,^{bi} Andres D. Perez,¹⁴ Marnel Snewe^{1,p} and Alejandro Szyukman^{1,...}

¹Bartinov de Fiskos Teletics U-DICONC, Vinolita Clamos, 1945; Garqueso Contalisano, 2004; Maleida, Spaine ²BTP, CONCET - Dytos de Fisico, Universidad Nacional de La Plano, C.C. 67, 1940 La Fisico, Algoritan ¹Departamente de Fisico, Teletica Fisico, Algoritan ²Departamente de Fisico, Teletica Control, Universida de Maleida (-2000) Cantobanco, Maleida (Spain ⁴Departamente de Pisico, Nicola (Clamos), ²Departamente de Pisico, Nicola (Clamos), ⁴Departamente de Pisico, Nicola (Clamos), ⁴Departemente de Pisico, Nicola (Clamos), ⁴

Imposing exclusion limits on new physics with machine-learned likelihoods

Ernesto Arganda,^{e,b} Martin de los Rios,^{e,c} Andres D. Perez^b and Rosa Maria Sandá Secane^{e,e}

⁴Instituto de Fisica Teórica UAH-CSIC, Cr/Neshis Cabren 13-15, Campus de Cantobhanco, 28049, Madrid, Spain Version Martin, Spain Martin,

^bIFLP, CONICET - Dyns. de Finica, Universidad Nacional de La Plata,

PS

About ML... Motivation ML Likelihood method Kernel Density Estimators

Traditional vs ML search of New Physics

Distinguish SM (bckg) vs BSM (signal) in collider data:

- Design observables, define control regions... \longrightarrow ML classifiers \checkmark
- For experimental significances, selection cuts \longrightarrow Working points X

э.

A 3 b

About ML... Motivation ML Likelihood method Kernel Density Estimators

Traditional vs ML search of New Physics

Distinguish SM (bckg) vs BSM (signal) in collider data:

- Design observables, define control regions... \longrightarrow ML classifiers \checkmark
- For experimental significances, selection cuts \longrightarrow Working points X



프 > 프

About ML... Motivation ML Likelihood method Kernel Density Estimators

Traditional vs ML search of New Physics

Distinguish SM (bckg) vs BSM (signal) in collider data:

- Design observables, define control regions... \longrightarrow ML classifiers \checkmark
- For experimental significances, selection cuts \longrightarrow Working points X

Is it possible to connect the ML classifier output with the standard statistical tests without defining working points?

→ Machine-Learned Likelihood (MLL) Method

About ML... Motivation ML Likelihood method Kernel Density Estimators

Traditional vs ML search of New Physics

Distinguish SM (bckg) vs BSM (signal) in collider data:

- Design observables, define control regions... \longrightarrow ML classifiers \checkmark
- For experimental significances, selection cuts \longrightarrow Working points X

Is it possible to connect the ML classifier output with the standard statistical tests without defining working points?

→ Machine-Learned Likelihood (MLL) Method

Can we avoid binning the output?

 \rightarrow +Kernel Density Estimators (KDE)

About ML... Motivation ML Likelihood method Kernel Density Estimators

The MLL method

Statistical model for ${\it N}$ independent measurements, with a high-dimensional set of observables x

$$\mathcal{L}(\mu, s, b) = p(N, \{x_i, i = 1, ..., N\} | \mu, s, b) \equiv \mathsf{Poiss}(N | \mu S + B) \prod_{i=1}^{N} p(x_i | \mu, s, b)$$

where S(B) is the expected total signal (background) yield, and

$$p(x|\mu, s, b) = \frac{B}{\mu S + B} p_b(x) + \frac{\mu S}{\mu S + B} p_s(x)$$

The relevant to derive exclusion limits on μ (considering models with $\mu \ge 0$)

$$\tilde{q}_{\mu} = \begin{cases} 0 & \text{if } \hat{\mu} > \mu \text{,} \\ -2 \text{ Ln } \frac{\mathcal{L}(\mu, s, b)}{\mathcal{L}(\hat{\mu}, s, b)} & \text{if } 0 \leqslant \hat{\mu} \leqslant \mu \text{,} \\ -2 \text{ Ln } \frac{\mathcal{L}(\mu, s, b)}{\mathcal{L}(0, s, b)} & \text{if } \hat{\mu} < 0 \text{,} \end{cases}$$

where $\hat{\mu}$ is the parameter that maximizes the likelihood

∃ >

About ML... Motivation ML Likelihood method Kernel Density Estimators

The MLL method

Statistical model for ${\it N}$ independent measurements, with a high-dimensional set of observables x

$$\mathcal{L}(\mu, s, b) = p(N, \{x_i, i = 1, ..., N\} | \mu, s, b) \equiv \mathsf{Poiss}(N|\mu S + B) \prod_{i=1}^{N} p(x_i | \mu, s, b)$$

where S(B) is the expected total signal (background) yield, and

$$p(x|\mu, s, b) = \frac{B}{\mu S + B} p_b(x) + \frac{\mu S}{\mu S + B} p_s(x)$$

The relevant to derive exclusion limits on μ (considering models with $\mu \ge 0$)

$$\tilde{q}_{\mu} = \begin{cases} 0 & \text{if } \hat{\mu} > \mu \\ 2(\mu - \hat{\mu})S - 2\sum_{i=1}^{N} \operatorname{Ln}\left(\frac{Bp_{b}(x_{i}) + \mu Sp_{s}(x_{i})}{Bp_{b}(x_{i}) + \hat{\mu}Sp_{s}(x_{i})}\right) & \text{if } 0 \leqslant \hat{\mu} \leqslant \mu \\ 2\mu S - 2\sum_{i=1}^{N} \operatorname{Ln}\left(1 + \frac{\mu Sp_{s}(x_{i})}{Bp_{b}(x_{i})}\right) & \text{if } \hat{\mu} < 0; \end{cases}$$

where $\hat{\mu}$ is the parameter that maximizes the likelihood

$$\sum_{i=1}^{N} \frac{p_{s}(x_{i})}{\hat{\mu} S \, p_{s}(x_{i}) + B \, p_{b}(x_{i})} = 1$$

4 E b

About ML... Motivation ML Likelihood method Kernel Density Estimators

The MLL method

Statistical model for ${\it N}$ independent measurements, with a high-dimensional set of observables x

$$\mathcal{L}(\mu, s, b) = p(N, \{x_i, i = 1, ..., N\} | \mu, s, b) \equiv \mathsf{Poiss}(N|\mu S + B) \prod_{i=1}^{N} p(x_i | \mu, s, b)$$

where S(B) is the expected total signal (background) yield, and

$$p(x|\mu, s, b) = \frac{B}{\mu S + B} p_b(x) + \frac{\mu S}{\mu S + B} p_s(x)$$

The relevant to derive exclusion limits on μ (considering models with $\mu \ge 0$)

$$\tilde{q}_{\mu} = \begin{cases} 0 & \text{if } \hat{\mu} > \mu \\ 2(\mu - \hat{\mu})S - 2\sum_{i=1}^{N} \operatorname{Ln}\left(\frac{Bp_{b}(x_{i}) + \mu Sp_{s}(x_{i})}{Bp_{b}(x_{i}) + \hat{\mu}Sp_{s}(x_{i})}\right) & \text{if } 0 \leqslant \hat{\mu} \leqslant \mu \\ 2\mu S - 2\sum_{i=1}^{N} \operatorname{Ln}\left(1 + \frac{\mu Sp_{s}(x_{i})}{Bp_{b}(x_{i})}\right) & \text{if } \hat{\mu} < 0; \end{cases}$$

where $\hat{\mu}$ is the parameter that maximizes the likelihood

$$\sum_{i=1}^{N} \frac{p_{s}(x_{i})}{\hat{\mu}S\,p_{s}(x_{i})+B\,p_{b}(x_{i})} = 1$$

4 E b

Introduction About ML... Applications ML Likelihood method Conclusions Kernel Density Estimator:

Solution: train classifier to distinguish signal from bckg with a balanced dataset. The classification score maximizing the binary cross-entropy and thus approaches

$$o(x) = \frac{p_s(x)}{p_s(x) + p_b(x)}$$

Dimensional reduction by dealing with o(x) instead of x

 $p_{s}(x) \rightarrow \tilde{p}_{s}(o(x))$, and $p_{b}(x) \rightarrow \tilde{p}_{b}(o(x))$





where $\tilde{p}_{s,b}(o(x))$ are the distributions of o(x) for signal and background, obtained by evaluating the classifier on a set of pure signal or background events, respectively.

Introduction Applications Conclusions	About ML Motivation
	ML Likelihood method
	Kernel Density Estimators

The relevant test statistic for exclusion limits

$$\tilde{q}_{\mu} = \begin{cases} 0 & \text{if } \hat{\mu} > \mu \\ 2(\mu - \hat{\mu})S - 2\sum_{i=1}^{N} \operatorname{Ln}\left(\frac{B\bar{p}_{b}(o(\mathbf{x}_{i})) + \mu S\bar{p}_{s}(o(\mathbf{x}_{i}))}{B\bar{p}_{b}(o(\mathbf{x}_{i})) + \hat{\mu}S\bar{p}_{s}(o(\mathbf{x}_{i}))}\right) & \text{if } 0 \leqslant \hat{\mu} \leqslant \mu \\ 2\mu S - 2\sum_{i=1}^{N} \operatorname{Ln}\left(1 + \frac{\mu S\bar{p}_{s}(o(\mathbf{x}_{i}))}{B\bar{p}_{b}(o(\mathbf{x}_{i}))}\right) & \text{if } \hat{\mu} < 0; \end{cases}$$

with $\hat{\mu}$ such us

$$\sum_{i=1}^{N} \frac{\tilde{p}_{s}(o(x_{i}))}{\hat{\mu}S\,\tilde{p}_{s}(o(x_{i})) + B\,\tilde{p}_{b}(o(x_{i}))} = 1$$

The median expected exclusion significance when the true hyphothesis is assumed to be the bckg-only one ($\mu^\prime=$ 0) is

$$\mathsf{med}\,\,[Z_\mu|\mathsf{0}] = \sqrt{\mathsf{med}\,\,[\tilde{q}_\mu|\mathsf{0}]}$$

About ML... Motivation ML Likelihood method Kernel Density Estimators

Traditional Binned-Likelihood (BL) method

 $p_{s,b}(x)/\tilde{p}_{s,b}(o(x_i))$ are not known and are approximated by discrete binned distributions D

$$\mathcal{L}(\mu, s, b) = \prod_{d=1}^{b} \mathsf{Poiss}(N_d | \mu S_d + B_d)$$

The median exclusion significance using Asimov datasets is given by



About ML... Motivation ML Likelihood method Kernel Density Estimators

What is the best way to extract $\tilde{p}_{s}(o(x))$ and $\tilde{p}_{b}(o(x))$?



æ

-

About ML... Motivation ML Likelihood method Kernel Density Estimators

 \longrightarrow Density estimation in a sense is the reverse of sampling: from given samples we want to retrieve the density function from which the samples were generated.

 \longrightarrow Two types of methods for density estimation

- Parametric: model the density function as a specified functional form with a fixed number of tunable parameters.
- Non-parametric: specify a model whose complexity grows with the number of training datapoints.

troduction pplications Conclusions	About ML Motivation
	ML Likelihood method
	Remer Density Estimate

Kernel Density Estimators

Kernel Density Estimators (KDE) is a non-parametric method for extracting $\tilde{p}_{s}(o(x_{i}))$ and $\tilde{p}_{b}(o(x_{i}))$

 \longrightarrow Smoothed version of the empirical distribution $q_o(x)$ of the training data $\{x_i, i = 1, ..., N\}$

$$q_o(x) = \frac{1}{N} \sum_{i}^{N} \delta(x - x_i)$$

 \longrightarrow We can smooth out the empirical distribution and turn it into a density by replacing each delta distribution with a smoothing kernel

$$\kappa_{\epsilon}(u) = \frac{1}{\epsilon^{D}} \kappa_{1}\left(\frac{u}{\epsilon}\right)$$

where $\epsilon > 0$ (bandwidth parameter) controls the width of the kernel and $\kappa_1(u)$ is a density function bounded from above (as $\epsilon \to 0$, $\kappa_{\epsilon}(u)$ approaches $\delta(u)$)

$$q_{\epsilon}(x) = \frac{1}{N} \sum_{i}^{N} \kappa_{\epsilon} (x - x_{i})$$

About ML... Motivation ML Likelihood method Kernel Density Estimators

$$\tilde{\rho}_{s,b}(o(x)) = \frac{1}{N} \sum_{i}^{N} \kappa_{\epsilon} \left[o(x) - o(x_i) \right]$$

Several options for κ_{ϵ} , e.g.



Introduction Applications Conclusions	About ML Motivation ML Likelihood method
	Kernel Density Estimators

The width parameter ϵ controls the degree of smoothness, if ϵ is too low the model may overfit, whereas if ϵ is too high the model may underfit. In general, we want ϵ to be smaller the more data we have and larger the higher the dimension is.



Toy Model: Multivariate Gaussian distributions Search for a heavy SSM Z' in dilepton final states at the HL-LHC

Toy Model: Multivariate Gaussian distributions

- Toy model in abstract space (x_1, x_2) . Events generated by $\mathcal{N}_2(\boldsymbol{m}, \boldsymbol{\Sigma})$ (known generative functions $p_{s,b}(x)$).
- Covariance matrices $\Sigma = \mathbb{I}_{2 \times 2}$ (no correlation) and m = +0.3(-0.3) $\mathbb{1}_2$ for *S* (*B*).



• Training of supervised per-event classifier, XGBoost with 1M events per class. KDE with Epanechnikov kernel.

Toy Model: Multivariate Gaussian distributions Search for a heavy SSM Z' in dilepton final states at the HL-LHC



- Significances estimated with all method are all close to the true pdf scenario given the low dimensionality of the problem.
- The ML output is always 1D regardless the dimensionality of the data and can be easyly handled.

∃ >

For higher dimensional data of $\dim = n$ with n > 2, $\mathcal{N}_n(\boldsymbol{m}, \boldsymbol{\Sigma})$, $\boldsymbol{\Sigma} = \mathbb{I}_{n \times n}$, $\boldsymbol{m} = \{+0.3(-0.3), +0.7(-0.7)\} \mathbb{1}_n$ for S(B):



- Results with MLL+KDE method approach the ones with the true generative functions, independently of dimension and/or separation of *S* and *B*.
- BL intractable in the original space. BL in the ML output not as good as MLL+KDE. Undesirable dependence with the binning.

Search for a heavy SSM Z' in dilepton final states at the HL-LHC

Based on ATLAS projections at $\sqrt{s} = 14$ TeV and 3000 fb⁻¹ for 95% CL exclusion limits on a Z'_{SSM} (ATLAS-PHYS-PUB-2018-014).

S:
$$p p \to Z' \to \ell^+ \ell^-$$

B: $p p \to Z/\gamma^* \to \ell^+ \ell^-$

- Use of $|p_T|$, ϕ , and η of the final state leptons as ML inputs in each channel.
- Training of supervised per-event classifier, XGBoost with 1M events per class. KDE with Epanechnikov kernel.
- For saving computational resources, generation of events only with dilepton invariant mass over 1.8 TeV.



Toy Model: Multivariate Gaussian distributions Search for a heavy SSM Z' in dilepton final states at the HL-LHC



- In both channels, unbinning signal and background posteriors provide more constraining limits than binning output.
- For direct comparison with ATLAS projections, necessary to simulate full spectrum of invariant masses. Potentially biased results by possible enhancement of classifier performance.

Conclusions

- MLL method allows to obtain exclusion (and discovery) significances for additive new physics scenarios.
- Uses a single XGBoost classifier and its full 1D output (no working points), which allows the estimation of the *S* and *B* pdfs needed for statistical inference. Not strictly necessary to bin the output to extract the pdfs.
- Inclusion of KDE as extension of MLL method to avoid the binning of the ML classifier output.
- Improves results obtained by traditional techniques in toy models and realistic analysis, approaching (when possible) the ones computed with true generative functions.
- Possible improvements: unsupervised analysis, systematic uncertainties...

Thank you!



(special thanks to M. de los Rios and A. D. Perez)

æ