hep-ml Tilman Plehn

- Losses Generati Control
- Testing
- SRegression



Tilman Plehn

Universität Heidelberg

Freiburg RTG, October 2023



Testing

SRegression

Transformers

Extracting correlations

 $\cdot\,$ Start with (compact) query representation

$$x_i \longrightarrow q = \frac{x_i}{|x|}$$

· Orthonormal values basis [related to *q* through scalar product]

$$q = \sum_j (q \cdot v_j) v_j$$

· Simpler orthogonal keys basis

$$q = \sum_{j} (q \cdot k_j) v_j$$
 with $k_j = \frac{v_j}{v^2}$

 $\rightarrow\,$ Self-attention representation

$$x_i \longrightarrow z_i = \sum_j (q \cdot k_j) v_j$$



Testing

Transformers

Extracting correlations

 $\cdot\,$ Start with (compact) query representation

$$x_i \longrightarrow q = \frac{x_i}{|x|}$$

• Orthonormal values basis [related to *q* through scalar product]

$$q = \sum_j (q \cdot v_j) v_j$$

· Simpler orthogonal keys basis

$$q = \sum_{j} (q \cdot k_j) v_j$$
 with $k_j = \frac{v_j}{v^2}$

→ Self-attention representation

$$x_i \longrightarrow z_i = \sum_j (q \cdot k_j) v_j$$

LHC phase space

- · learn bin-bin relation $x_i \leftrightarrow x_j$
- · latent query representation $q = W^Q x$ latent key representation $k = W^K x$ correlation $A_{ij} = q_i \cdot k_j$
- · latent value representation $v = W^V x$ constructed representation z = A v





hep-ml Tilman Plehn JetGPT Losses Generation Control

JetGPT

Autoregressive transformer

· factorized density

$$p_{\text{model}}(x|\theta) = \prod_{i} p(x_i|x_1, ..., x_{i-1})$$

- $\cdot \ \text{bins} \rightarrow \text{Gaussian}$ mixture model
- · autoregressive $A_{ij} = 0$ for j > i
- → Bayesian version for uncertainties





hep-ml Tilman Plehn JetGPT Losses Generation Control

JetGPT

Autoregressive transformer

· factorized density

$$p_{\text{model}}(x|\theta) = \prod_{i} p(x_i|x_1,...,x_{i-1})$$

- $\cdot \ \text{bins} \rightarrow \text{Gaussian}$ mixture model
- · autoregressive $A_{ij} = 0$ for j > i
- → Bayesian version for uncertainties



Bayesian JetGPT

· sometimes you win...





hep-ml Tilman Plehn JetGPT Losses Generation

JetGPT

Autoregressive transformer

· factorized density

$$p_{\text{model}}(x|\theta) = \prod_{i} p(x_i|x_1,...,x_{i-1})$$

- $\cdot \ \text{bins} \rightarrow \text{Gaussian}$ mixture model
- · autoregressive $A_{ij} = 0$ for j > i
- → Bayesian version for uncertainties



Bayesian JetGPT

- · sometimes you win...
 - ...and sometimes there is work to do...





Testing

Likelihood loss & uncertainties

Loss to train θ -distributions

- · energy measurement for jet j $\langle E \rangle = \int dE \ E \ p(E)$
- · weighted by reproduced training data $p(\theta|T)$ $p(E) = \int d\theta \ p(E|\theta) \ p(\theta|T)$
- $\rightarrow \theta$ -distributions means Bayesian NN



Likelihood loss & uncertainties

Loss to train θ -distributions

- · energy measurement for jet j $\langle E \rangle = \int dE \ E \ p(E)$
- · weighted by reproduced training data $p(\theta|T)$ $p(E) = \int d\theta \ p(E|\theta) \ p(\theta|T)$
- $\rightarrow \theta$ -distributions means Bayesian NN

Variational approximation

definition of training [think
$$q(\theta)$$
 as Gaussian with mean and width]
 $p(E) = \int d\theta \ p(E|\theta) \ p(\theta|T) \approx \int d\theta \ p(E|\theta) \ q(\theta)$

 $\begin{array}{ll} \cdot \mbox{ similarity through minimal KL-divergence } & \mbox{[Bayes' theorem to remove unknown posterior]} \\ D_{\text{KL}}[q(\theta), p(\theta|T)] &= \int d\theta \ q(\theta) \ \log \frac{q(\theta)}{p(\theta|T)} \\ &= \int d\theta \ q(\theta) \ \log \frac{q(\theta)p(T)}{p(T|\theta)p(\theta)} \\ &= D_{\text{KL}}[q(\theta), p(\theta)] - \int d\theta \ q(\theta) \ \log p(T|\theta) + \log p(T) \int d\theta \ q(\theta) \end{array}$



Likelihood loss & uncertainties

Loss to train θ -distributions

- · energy measurement for jet j $\langle E \rangle = \int dE \ E \ p(E)$
- · weighted by reproduced training data $p(\theta|T)$ $p(E) = \int d\theta \ p(E|\theta) \ p(\theta|T)$
- $\rightarrow \theta$ -distributions means Bayesian NN

Variational approximation

- · definition of training [think $q(\theta)$ as Gaussian with mean and width] $p(E) = \int d\theta \ p(E|\theta) \ p(\theta|T) \approx \int d\theta \ p(E|\theta) \ q(\theta)$
- · similarity through minimal KL-divergence [Bayes' theorem to remove unknown posterior]

$$\begin{split} D_{\mathsf{KL}}[q(\theta), p(\theta|T)] &= \int d\theta \ q(\theta) \ \log \frac{q(\theta)}{p(\theta|T)} \\ &= \int d\theta \ q(\theta) \ \log \frac{q(\theta)p(T)}{p(T|\theta)p(\theta)} \\ &= D_{\mathsf{KL}}[q(\theta), p(\theta)] - \int d\theta \ q(\theta) \ \log p(T|\theta) + \log p(T) \end{split}$$



Likelihood loss & uncertainties

Loss to train θ -distributions

- · energy measurement for jet j $\langle E \rangle = \int dE \ E \ p(E)$
- · weighted by reproduced training data $p(\theta|T)$ $p(E) = \int d\theta \ p(E|\theta) \ p(\theta|T)$
- $\rightarrow \theta$ -distributions means Bayesian NN

Variational approximation

definition of training [think
$$q(\theta)$$
 as Gaussian with mean and width]
 $p(E) = \int d\theta \ p(E|\theta) \ p(\theta|T) \approx \int d\theta \ p(E|\theta) \ q(\theta)$

· similarity through minimal KL-divergence [Bayes' theorem to remove unknown posterior]

$$\begin{split} D_{\mathsf{KL}}[q(\theta), p(\theta|T)] &= \int d\theta \ q(\theta) \ \log \frac{q(\theta)}{p(\theta|T)} \\ &= \int d\theta \ q(\theta) \ \log \frac{q(\theta)p(T)}{p(T|\theta)p(\theta)} \\ &\approx D_{\mathsf{KL}}[q(\theta), p(\theta)] - \int d\theta \ q(\theta) \ \log p(T|\theta) \equiv \mathcal{L} \end{split}$$



→ Two-term loss: likelihood + prior

hep-ml Tilman Plehn JetGPT Losses Generation Control

Relation to deterministic networks

Regularization

· BNN loss

$$\mathcal{L} = -\int d heta \; q(heta) \; \log p(T| heta) + D_{\mathsf{KL}}[q(heta), p(heta)]$$



hep-ml Tilman Plehn JetGPT Losses Generation Control

Relation to deterministic networks

Regularization

· Gaussian prior

$$\mathcal{L} = -\int d\theta \; q(\theta) \; \log p(T|\theta) + rac{\sigma_q^2 - \sigma_p^2 + (\mu_q - \mu_p)^2}{2\sigma_p^2} + \dots$$

 \cdot deterministic network

$$q(heta) = \delta(heta - heta_0) \quad \Rightarrow \quad \mathcal{L} pprox - \log p(T| heta_0) + rac{(heta_0 - \mu_p)^2}{2\sigma_p^2}$$

 $\rightarrow\,$ Likelihood with L2-regularization



SRegression

Relation to deterministic networks

Regularization

· Gaussian prior

$$\mathcal{L} = -\int d\theta \ q(\theta) \ \log p(T|\theta) + \frac{\sigma_q^2 - \sigma_p^2 + (\mu_q - \mu_p)^2}{2\sigma_p^2} + \dots$$

· deterministic network

$$q(heta) = \delta(heta - heta_0) \quad \Rightarrow \quad \mathcal{L} pprox - \log p(T| heta_0) + rac{(heta_0 - \mu_p)^2}{2\sigma_p^2}$$

 \rightarrow Likelihood with L2-regularization

Dropout

· Bernoulli weights

$$q(\theta) \rightarrow q(x) = \rho^{x} (1-\rho)^{1-x} \bigg|_{x=0,1}$$
 with $\theta = x\theta_{0}$

· likelihood loss

$$\mathcal{L} = -\sum_{x=0,1} \left[\rho^x (1-\rho)^{1-x} \right] \log p(T|x\theta_0) = -\rho \log p(T|\theta_0)$$

- · likelihood Gaussian or whatever else...
- \rightarrow Regularized likelihood with dropout



hep-ml Tilman Plehn JetGPT Losses

Control Uncertainty Testing

CRograndia

Statistics vs systematics

Network evaluation

· expectation value using trained network $q(\theta)$

$$\langle E \rangle = \int dEd\theta \ E \ p(E|\theta) \ q(\theta)$$

$$\equiv \int d\theta \ q(\theta)\overline{E}(\theta) \quad \text{with} \quad \overline{E}(\theta) = \int dE \ E \ p(E|\theta)$$

· NN-variance

$$\begin{aligned} \sigma_{\text{tot}}^{2} &= \int dEd\theta \ \left(E - \langle E \rangle\right)^{2} \ p(E|\theta) \ q(\theta) \\ &= \int d\theta \ q(\theta) \left[\overline{E^{2}}(\theta) - 2\langle E \rangle \overline{E}(\theta) + \langle E \rangle^{2}\right] \\ &= \int d\theta \ q(\theta) \left[\overline{E^{2}}(\theta) - \overline{E}(\theta)^{2} + \left(\overline{E}(\theta) - \langle E \rangle\right)^{2}\right] \equiv \sigma_{\text{stoch}}^{2} + \sigma_{\text{pred}}^{2} \end{aligned}$$

Two uncertainties

· statistical — vanishing for $q(\theta) \rightarrow \delta(\theta - \theta_0)$

$$\sigma_{\mathsf{pred}}^2 = \int d\theta \; q(\theta) \left[\overline{E}(\theta) - \langle E \rangle\right]^2$$

 \cdot systematic — vanishing for $p(E| heta)
ightarrow \delta(E-E_0)$

$$\sigma_{\text{stoch}}^2 \equiv \sigma_{\text{model}}^2 = \int d\theta \ q(\theta) \left[\overline{E^2}(\theta) - \overline{E}(\theta)^2\right]$$



resung

Jets and non-Gaussian errors

Measure $p_{T,t}$ of hadronically decaying top

- BNN regression $p_{T,t}$ p_T of (fat) jet decent estimate for $p_{T,t}^{\text{truth}}$
- non-Gaussian truth label

symmetric in ISR-jet 'QCD heat bath' without ISR jets need for correction





SRearessior

Jets and non-Gaussian errors

Measure $p_{T,t}$ of hadronically decaying top

- BNN regression $p_{T,t}$ p_T of (fat) jet decent estimate for $p_{T,t}^{\text{truth}}$
- non-Gaussian truth label symmetric in ISR-jet 'QCD heat bath' without ISR jets need for correction
- · training sample size

 $\begin{array}{l} \text{separate } \sigma_{\text{stoch}} \gg \sigma_{\text{pred}} \\ \text{statistics not the problem} & \text{[LHC theme]} \\ \text{noisy label inherent limitation} \\ \text{checked with ensembles} \end{array}$





SRegression

Jets and non-Gaussian errors

Measure $p_{T,t}$ of hadronically decaying top

- BNN regression $p_{T,t}$ p_T of (fat) jet decent estimate for $p_{T,t}^{\text{truth}}$
- · non-Gaussian truth label

symmetric in ISR-jet 'QCD heat bath' without ISR jets need for correction

· training sample size

 $\begin{array}{l} \text{separate } \sigma_{\text{stoch}} \gg \sigma_{\text{pred}} \\ \text{statistics not the problem} & \text{[LHC theme]} \\ \text{noisy label inherent limitation} \\ \text{checked with ensembles} \end{array}$

non-Gaussian network output

remember $p_{T,t}^{\text{truth}}$ non-Gaussian model $p(T|\theta)$ as Gaussian mixture weight distribution $q(\theta)$ still Gaussian





Generative networks

Unsupervised Bayesian networks

- data: event sample [points in 2D space]
 learn phase space density normalizing flow mapping to latent space [INN] standard distribution in latent space [Gaussian] mapping bijective sample from latent space
- Bayesian version allow weight distributions learn uncertainty map
- · 2D wedge ramp

$$p(x) = ax + b = ax + \frac{1 - \frac{a}{2}(x_{\max}^2 - x_{\min}^2)}{x_{\max} - x_{\min}}$$
$$(\Delta p)^2 = \left(x - \frac{1}{2}\right)^2 (\Delta a)^2 + \left(1 + \frac{a}{2}\right)^2 (\Delta x_{\max})^2 + \left(1 - \frac{a}{2}\right)^2 (\Delta x_{\min})^2$$

explaining minimum in $\sigma_{\text{pred}}(x)$





 \rightarrow INNs just (non-parametric) fits



More generative networks

Alternative architectures

- · always a fit?
- · expressivity vs architecture?
- $\cdot\,$ discrete diffusion model





More generative networks

Alternative architectures

- · always a fit?
- · expressivity vs architecture?
- · discrete diffusion model
- · continuous diffusion model





SRegression

More generative networks

Alternative architectures

- · always a fit?
- · expressivity vs architecture?
- · discrete diffusion model
- · continuous diffusion model
- · autoregressive transformer (bins)





Precision generator

Phase-space generators [typical LHC task]

- training from event samples no energy-momentum conservation
- \cdot every correlation counts
- $\cdot ~Z_{\mu\mu} + \{1,2,3\}~ ext{jets}~$ [Z-peak, variable jet number, jet-jet topology]



Precision generator

Phase-space generators [typical LHC task]

- training from event samples no energy-momentum conservation
- $\cdot\,$ every correlation counts
- $\cdot \,\, Z_{\mu\mu} + \{1,2,3\} \,\, ext{jets} \,\,$ [Z-peak, variable jet number, jet-jet topology]

INN-generator

stable bijective mapping

 $G_{\theta}(r) \rightarrow$ latent $r \sim p_{\text{latent}}$ phase space $x \sim p_{data}$ $\leftarrow \overline{G}_{\theta}(x)$ tractable Jacobian Z + 1 jet exclusive nor 10-5 10-5 $dx p_{model}(x) = dr p_{latent}(r)$ Truth $\frac{\partial \overline{G}_{\theta}(x)}{\partial x}$ $p_{\text{model}}(x) = p_{\text{latent}}(\overline{G}_{\theta}(x))$ INN Train 10^{-4} likelihood loss $\mathcal{L}_{\text{INN}} = -\Big\langle \log p_{\text{model}}(x) \Big\rangle_{p_{\text{data}}}$ the state of the s

10.0 1.0

0.1

50

125

150

100

 p_{T,j_1} [GeV]

 \Rightarrow Per-cent precision possible



hep-ml Tilman Plehn JetGPT Losses Generation Control

Uncertai

- Testing
- SRegressio

Control and reweighting

Best of GANs: discriminator

- · input $\{p_T, \eta, \phi, M, M_{\mu\mu}, \Delta R\}$
- · output D = 0 (generator) vs D = 1 (training)
- · NP-optimal discriminator

$$D(x)
ightarrow rac{p_{ ext{data}}(x)}{p_{ ext{data}}(x) + p_{ ext{model}}(x)}
ightarrow rac{1}{2}$$

· learned event weight

$$w(x)
ightarrow rac{D(x)}{1-D(x)} = rac{
ho_{ ext{data}}(x)}{
ho_{ ext{model}}(x)}
ightarrow 1$$

 \Rightarrow Dual purpose: control and reweight







lesting

SRegression

Events with uncertainties

Bayesian network generator

- network with weight distributions [Gal sample weights [defining error bar] frequentist: efficient ensembling
- \Rightarrow Training-related error bars





hep-ml Tilman Plet JetGPT Losses Generation Control Uncertainty Testing SRegression

Events with uncertainties

Bayesian network generator

- network with weight distributions [Gal (2016)] sample weights [defining error bar] frequentist: efficient ensembling
- \Rightarrow Training-related error bars

Theory uncertainties

- BNN regression/classification: systematics from data augmentation
- · systematic uncertainties in tails

$$w = 1 + a \left(\frac{p_{T,j_1} - 15 \text{ GeV}}{100 \text{ GeV}}\right)^2$$

- augment training data $[a = 0 \dots 30]$
- train conditionally on a error bar from sampling a
- ⇒ Systematic/theory error bars





Events with uncertainties

Bayesian network generator

- network with weight distributions [Ga sample weights [defining error bar] frequentist: efficient ensembling
- \Rightarrow Training-related error bars

Theory uncertainties

- BNN regression/classification: systematics from data augmentation
- · systematic uncertainties in tails

$$w = 1 + a \left(\frac{p_{T,j_1} - 15 \text{ GeV}}{100 \text{ GeV}}\right)^2$$

- augment training data $[a = 0 \dots 30]$
- train conditionally on a error bar from sampling a
- ⇒ Systematic/theory error bars





Testing generative networks

Compare network to training/test data

- · supervised: histogram deviation [or pull]
- $\cdot \,$ unsupervised density \rightarrow histogram discriminator

$$w(x_i) = \frac{D(x_i)}{1 - D(x_i)} = \frac{p_{\text{data}}(x_i)}{p_{\text{model}}(x_i)}$$

 $\rightarrow\,$ Using interpretable phase space



Testing generative networks

Compare network to training/test data

- · supervised: histogram deviation [or pull]
- $\cdot \,$ unsupervised density \rightarrow histogram discriminator

$$w(x_i) = \frac{D(x_i)}{1 - D(x_i)} = \frac{p_{\text{data}}(x_i)}{p_{\text{model}}(x_i)}$$

 \rightarrow Using interpretable phase space

Applied to event generators [also jets, calorimeter showers]

- · shape and width of w-histogram
- · pattern in (interpretable) phase space?







Testing generative networks

Compare network to training/test data

- · supervised: histogram deviation [or pull]
- $\cdot \,$ unsupervised density \rightarrow histogram discriminator

$$w(x_i) = \frac{D(x_i)}{1 - D(x_i)} = \frac{p_{\text{data}}(x_i)}{p_{\text{model}}(x_i)}$$

 \rightarrow Using interpretable phase space

Applied to event generators [also jets, calorimeter showers]

- · shape and width of w-histogram
- · pattern in (interpretable) phase space?







 \rightarrow Generative xAI for LHC physicists

hep-ml Tilman Plehn JetGPT Losses Generation

- Uncortai
- Testing

Symbolic regression

Symbolic regression of score [PySR (M Cranmer) + final fit]

- · function to approximate $t(x|\theta)$
- \cdot phase space parameters $x_p = p_T/m_H, \Delta\eta, \Delta\phi$ [node]
- \cdot operators $\sin x, x^2, x^3, x + y, x y, x * y, x/y$ [node]
- · represent formula as tree [complexity = number of nodes]
- \Rightarrow Figures of merit

$$\mathsf{MSE} = \frac{1}{n} \sum_{i=1}^{n} \left[g_i(x) - t(x, z|\theta) \right]^2 \to \mathsf{MSE} + \mathsf{parsimony} \cdot \mathsf{complexity}$$

Score around Standard Model

compl	dof	function	MSE	•
3	1	$a \Delta \phi$	$1.30\cdot 10^{-1}$	Λ.
4	1	$\sin(a\Delta\phi)$	$2.75\cdot 10^{-1}$	
5	1	$a\Delta\phi x_{p,1}$	$9.93 \cdot 10^{-2}$	10-1
6	1	$-x_{p,1}\sin(\Delta\phi+a)$	$1.90\cdot10^{-1}$	ш і 🔓
7	1	$(-x_{p,1}-a)\sin(\sin(\Delta\phi))$	$5.63 \cdot 10^{-2}$	W
8	1	$(a - x_{p,1})x_{p,2}\sin(\Delta\phi)$	$1.61 \cdot 10^{-2}$	
14	2	$x_{p,1}(a\Delta\phi - \sin(\sin(\Delta\phi)))(x_{p,2} + b)$	$1.44\cdot10^{-2}$	
15	3	$-(x_{p,2}(a\Delta\eta^2 + x_{p,1}) + b)\sin(\Delta\phi + c)$	$1.30 \cdot 10^{-2}$	· · · · ·
16	4	$-x_{p,1}(a-b\Delta\eta)(x_{p,2}+c)\sin(\Delta\phi+d)$	$8.50 \cdot 10^{-3}$	10-2
28	7	$\begin{vmatrix} (x_{p,2}+a)(bx_{p,1}(c-\Delta\phi) \\ -x_{p,1}(d\Delta\eta + ex_{p,2} + f)\sin(\Delta\phi + g)) \end{vmatrix}$	$8.18\cdot 10^{-3}$	5 10 15 20 complexity

25



hep-ml Tilman Pleh JetGPT Losses Generation Control

Testing

SRegression

Symbolic regression

Symbolic regression of score [PySR (M Cranmer) + final fit]

- function to approximate $t(x|\theta)$
- \cdot phase space parameters $x_{p}=p_{T}/m_{H},\Delta\eta,\Delta\phi$ [node]
- \cdot operators $\sin x, x^2, x^3, x + y, x y, x * y, x/y$ [node]
- · represent formula as tree [complexity = number of nodes]
- \Rightarrow Figures of merit

$$\mathsf{MSE} = \frac{1}{n} \sum_{i=1}^{n} \left[g_i(x) - t(x, z | \theta) \right]^2 \rightarrow \mathsf{MSE} + \mathsf{parsimony} \cdot \mathsf{complexity}$$

Score around Standard Model

- · expected limits:
 - very wrong formula wrong formula
- same within statistical limitation: right formula MadMiner
- ⇒ Formulas to numerics and back





hep-ml Tilman Plehn JetGPT Losses Generation

Uncertai

resting

SRegression

ML Uncertainties

ML-applications

- · just another numerical tool for a numerical field
- $\cdot\,$ driven by money from data science and medical research
- · goals are...

...improve established tasks ...develop new tools for established tasks ...transform through new ideas

- · xAI through...
 - ...precision control
 - ...uncertainties
 - ...symmetries
 - ...formulas

 $\rightarrow~$ Lots of fun with hard LHC problems

Modern Machine Learning for LHC Physicists

Tilman Plehna; Anja Buttera, Barry Dillona, Claudius Krausea, and Ramon Winterhalderd

^a Institut für Theoretische Physik, Universität Heidelberg, Germany ^b LPNHE, Sorbonne Université, Université Paris Cité, CNRS/IN2P3, Paris, France ^c NHETC, Dept. of Physics and Astronomy, Rutgers University, Piscataway, USA ^d CP3, Université Catholique de Louvain, Louvain-La-Neuve, Belgium

July 21, 2023

Abstract

Moders machine tearing in transforming particle physics, faster than we can fields, and bulyeng its way ium our merrical tool bics. For your generactives its its calls to any one op of also development, which means applying entitives edge methods, and tools to the full arrange of LHC physics problems. These lecture noise, are meant to load attents with popular. They are intra-field and the state of the physics. They are intra-field characterized problems in the state methods and to state of the distances are well and the local exciton state and the state and uncertainty states metrics. As part of the physics are stated and the state of the state of the state of the state of the distances and the state of the distances are stated and the state of the distances are stated and the state of the distance performance and the state of the distance of th

