Statistics for HEP

Part 2: Machine Learning

Invited lectures, 12th Course on Physics of the LHC (LIP, Lisboa, Portugal)

Dr. Pietro Vischia pietro.vischia@cern.ch @pietrovischia



If you are reading this as a web page: have fun! If you are reading this as a PDF: please visit

https://www.hep.uniovi.es/vischia/persistent/2024-03-20-

22_StatisticsAtCoursePhysicsLHC_vischia_part2.html

to get the version with working animations

Lecture 2 Artificial Intelligence

Pietro Vischia - Statistics for HEP (13th Course on Physics of the LHC, Lisboa, Portugal) - 2024.03.20-22 --- 2 / 108

Practicalities

- Significantly restructured with respect to the past years
 - Lecture 1: Probability and Statistics (minus hypothesis testing)
 - Lecture 2: Machine Learning (and hypothesis testing)
- More detailed material in my twenty-hours intensive course
 - It may be useful if you tried out the exercises, at your pace!
- Many references here and there, and in the last slide
 - Try to read some of the referenced papers!
 - Unreferenced stuff copyrighted P. Vischia for inclusion in my (finally) upcoming textbook

AI, the eternal buzzword?

- Artificial Intelligence (AI)
- Machine Learning (ML)
- Statistical Learning



Complex Experiments



Complex Data



Likelihood and information

• Data sample X_{obs}

$$\mathcal{L}(X; heta):=P(X| heta)|_{X_{obs}}$$

- The Likelihood Principle: The likelihood function $L(\vec{x}; \theta)$ contains all the information available in the data sample relevant for the estimation of θ
 - ✓ Bayesian statistics
 - X Frequentist statistics

$$I(heta) = -Eigg[\Big(rac{\partial^2}{\partial heta^2} ln L(X; heta) \Big)^2 | heta_{true} igg]$$



We like low-dim summaries



• Physical observable for inference



Brain activity...



...approximated...



...using computers



Santiago Ramón y Cajal



Santiago Ramón y Cajal











$$I = C rac{dV}{dt} + G_{Na} m^3 h (V - V_{Na}) + G_K n^4 (V - V_K) + G_L (V - V_L)$$

Computationally heavy



Simplified Neurons

Bulletin of Mathematical Biology Vol. 52, No. 1/2, pp. 99-115, 1990. Printed in Great Britain. 0092-8240/90\$3.00 + 0.00 Pergamon Press plc Society for Mathematical Biology

A LOGICAL CALCULUS OF THE IDEAS IMMANENT IN NERVOUS ACTIVITY*

WARREN S. MCCULLOCH AND WALTER PITTS University of Illinois, College of Medicine, Department of Psychiatry at the Illinois Neuropsychiatric Institute, University of Chicago, Chicago, U.S.A.

Because of the "all-or-none" character of nervous activity, neural events and the relations among them can be treated by means of propositional logic. It is found that the behavior of every net can be described in these terms, with the addition of more complicated logical means for nets containing circles; and that for any logical expression satisfying certain conditions, one can find a net behaving in the fashion it describes. It is shown that many particular choices among possible neurophysiological assumptions are equivalent, in the sense that for every net behaving under one assumption, there exists another net which behaves under the other and gives the same results, although perhaps not in the same time. Various applications of the calculus are discussed.

Perceptrons

$$y=f\Big(b_i+\sum w_ix_i\Big)$$



The perceptron

• Linear combination of the inputs

$$\sum_j w_j x_j$$

• Activation function (imitating activation potentials in cells)

$$g\Big(w_0+\sum_j w_j x_j\Big)$$

• Bias term (order zero in the inputs \rightarrow traslation)

$$\hat{y} = g\Big(w_0 + \sum_j w_j x_j\Big)$$

• In matrix form, it's a line-by-column product

$$\hat{y} = g\Big(w_0 + \mathbf{X}^T \mathbf{W}\Big)$$

Activation function

• The original activation function was linear

$$egin{aligned} g(z) &= 1 ext{ if } z = w_0 + \sum w_i x_i > = 0 \ g(z) &= -1 ext{ if } z = w_0 + \sum w_i x_i < 0 \end{aligned}$$

- Not possible to get nonlinear separation surfaces
- Sigmoid

$$g(z)=rac{1}{1+e^{-z}}$$
 , with derivative $g'(z)=g(z)\Big(1-g(z)\Big)$

• Rectified Linear Unit (ReLU)

g(z)=max(0,z) , with derivative $g^{\prime}(z)=1$ if z>0,0 otherwise

Artificial Neural Networks



Learn in different ways



Supervised learning



Gradient Descent

- Optimize/learn by finding the minimum of a function $\mathcal{L}:\mathbb{R}^n
 ightarrow\mathbb{R}$
- Nonconvex problems: saddle points, manifolds of minima
- Empirical risk minimization

 $\hat{L}(f) = rac{1}{n} \sum_{i=1}^n |f(x_i) - f^*(x_i)|^2$

• Generalization (for learning problems)





Vapnik's image from youtube, Gradient descent animation by Alec Radford

Organizing the space

- Formally, act on a space of function
- Need a notion of complexity.
 - $\circ \ \mathsf{Norm} \to \mathsf{Banach} \, \mathsf{space}!$



The Fundamental Theorem of Machine Learning

- We want to relate the result of the empirical risk minimization (ERM) with the prediction
 - Let's use the constraint form
- Let's assume we have solved the ERM at a precision ϵ (we are ϵ -away from...). We then have $\hat{f} \in \mathcal{F}^{\delta}$ such that $\hat{L}(\hat{f}) \leq \epsilon + \min_{f \in \mathcal{F}^{\delta}} \hat{L}(f)$
- How good is \hat{f} at predicting f^* ? In other words, what's the true loss?
 - Can use the triangular inequality

The Fundamental Theorem of Machine Learning

 $L(\widehat{f}) - \inf_{f \in \mathcal{F}} L(f) \leq \inf_{f \in \mathcal{F}^{\delta}} L(f) - \inf_{f \in \mathcal{F}} L(f)$

$$+2 \sup_{f\in \mathcal{F}^{\delta}} \left| L(f) - \hat{L}(f)
ight|$$

 $+\epsilon$

Approximation error

(how appropriate is my measure of complexity)

Statistical error

(impact of having the) empirical loss instead of the true loss)

Optimization error

Loss function

- Empirical Loss Function
- Cost function
- Empirical risk

$$\mathbf{J}(\mathbf{W}) = rac{1}{n}\sum_{i=1}^n \mathcal{L}(f(x^{(i)};\mathbf{W}),y^{*(i)})$$

• Minimized by:

$$\mathbf{W}^0 = argmin_\mathcal{W} \mathbf{J}(\mathbf{W}) = argmin_\mathcal{W} \mathbf{J}(\mathbf{W}) = rac{1}{n} \sum_{i=1}^n \mathcal{L}(f(x^{(i)};\mathbf{W}),y^{*(i)})$$

Backpropagation

$$\mathbf{J}(\mathbf{W}) = rac{1}{n}\sum_{i=1}^n \mathcal{L}(f(x^{(i)};\mathbf{W}),y^{*(i)}), \qquad \mathbf{W}^0 = argmin_{\mathcal{W}}\mathbf{J}(\mathbf{W})$$

 $\mathbf{W} \leftarrow \mathbf{W} + \eta rac{\partial \mathbf{J}(\mathbf{W})}{\partial \mathbf{W}}$



Derive



Automatic differentiation

has many names

- Automatic differentiation
- Algorithmic differentiation
- AD
- Autodiff
- Algodiff
- Autograd

Automatic differentiation

 $z(x,y) = 2x + x \sin(y) + y^3$



Forward mode

- To the extreme, $f:\mathbb{R}
 ightarrow\mathbb{R}^m$
- Evaluates $\left(\frac{\partial f_1}{\partial x}, \ldots, \frac{\partial f_m}{\partial x}\right)$

Reverse mode

- To the extreme, $f:\mathbb{R}^n
 ightarrow\mathbb{R}$
- Evaluate $abla f(\mathbf{x})(rac{\partial f}{\partial x_1},\ldots,rac{\partial f}{\partial x_n})$
- Computational cost of calculating $\mathbf{J}_f(\mathbf{x})$ for $f:\mathbb{R}^n o\mathbb{R}^m$ in $\mathbb{R}^n imes\mathbb{R}^m$

 $\mathcal{O}(n \, \operatorname{time}(f))$

 $\mathcal{O}(m \operatorname{time}(f))$

Forward and reverse (==backprop) modes

Primal: independent to dependent

Adjoint (derivatives): dependent to independent

$$y({f x})=2x_0+x_0\,sin(x_1)+x_1^3$$

<i>Fwd Primal Trace</i> Atomic operation	Value in $(1,2)$	Fwd Tangent Trace (set $\dot{x}_0 =$ 1 to compute $\frac{\partial y}{\partial x_0}$) Atomic operation			operation		operation	
			Value in $(1,2)$		$egin{array}{l} v_0 = x_0 \ v_1 = x_1 \end{array}$	$1 \\ 2$	$ar{x}_0 = ar{v}_0 \ ar{x}_1 = ar{v}_1$	$2.9093 \\ 11.5839$
					$egin{aligned} &v_2 = 2v_0 \ v_3 = \ sin(v_1) \ v_4 = \ v_0v_3 \ v_5 = v_1^3 \ v_6 = \ v_2 + \ v_4 + v_5 \end{aligned}$	$2 \\ 0.9093 \\ 0.9093 \\ 8 \\ 10.9093$	$ar{v}_0 = ar{v}_0 + ar{v}_2 \partial v_2 / \partial v_0 \ ar{v}_0 = ar{v}_4 \partial v_4 / \partial v_0 \ ar{v}_1 = ar{v}_1 + ar{v}_3 \partial v_3 / \partial v_1 \ ar{v}_1 = ar{v}_5 \partial v_5 / \partial v_1 \ ar{v}_2 = ar{v}_6 \partial v_6 / \partial v_2 \ ar{v}_3 = ar{v}_4 \partial v_4 / \partial v_3 \ ar{v}_4 = ar{v}_6 \partial v_6 / \partial v_4 \ ar{v}_5 = ar{v}_6 \partial v_6 / \partial v_5$	$egin{array}{lll} ar{v}_0 + ar{v}_2 imes \ 2 = \ 2.9093 \ ar{v}_4 imes v_3 = \end{array}$
$egin{array}{l} v_0 = x_0 \ v_1 = x_1 \end{array}$	$egin{array}{c} 1 \\ 2 \end{array}$	$egin{array}{lll} \dot{v}_0 = \dot{x}_0 \ \dot{v}_1 = \dot{x}_1 \end{array}$	1 0					$egin{array}{llllllllllllllllllllllllllllllllllll$
$egin{aligned} &v_2 = 2v_0 \ v_3 = \ sin(v_1) \ v_4 = \ v_0v_3 \ v_5 = v_1^3 \ v_6 = \ v_2 + \ v_4 + v_5 \end{aligned}$	$2 \\ 0.9093 \\ 0.9093 \\ 8 \\ 10.9093$	$egin{aligned} \dot{v}_2 &= 2\dot{v}_0 \ \dot{v}_3 &= \dot{v}_1 cos(v_1) \ \dot{v}_4 &= \dot{v}_0 v_3 + \ v_0 \dot{v}_3 \ \dot{v}_5 &= 3\dot{v}_1 v_1^2 \ \dot{v}_6 &= \dot{v}_2 + \dot{v}_4 + \ \dot{v}_5 \end{aligned}$	$\begin{array}{c} 2 \times 1 \\ 0 \times \\ -0.41 \\ 1 \times \\ 0.9093 + \\ 1 \times 0 \\ 3 \times 0 \times 4 \\ 2 + \\ 0.9093 + \\ 0 \end{array}$					
$y=v_6$	10.9093	$\dot{y}=\dot{v}_{6}$	2.9093				, -	1

Fwd

Primal

Trace

Atomic

Value in

(1, 2)

Pietro Vischia - Statigties f w_6 IEP (13 10.9093 n Ph \overline{y}_6 cs of \overline{y}_1 e LHC, Lisboa, Portugal 1-2024.03.20-22 --- 34 / 108

Rev Adjoint

Trace (set $\bar{y} =$

1 to compute

 $\frac{\partial v}{\partial y}$)

Value in

(1, 2)

Designed to be simple in software

```
import torch, math
x0 = torch.tensor(1., requires_grad=True)
x1 = torch.tensor(2., requires_grad=True)
p = 2*x0 + x0*torch.sin(x1) + x1**3
print(p)
p.backward()
print(x0.grad, x1.grad)
```

yielding

```
Primal: tensor(10.9093, grad_fn=<AddBackward0>)
Adjoint: tensor(2.9093) tensor(11.5839)
```

- Torch (and similar software) will correctly differentiate only when the atomic operations are supported within it
 - Common operations are overloaded (__mul__ rewritten by torch._mult_)
 - Operations from libraries (math.sin()) must be replaced by their differentiation-aware equivalents (torch.sin())
Learning as a tradeoff...

• ``a model with zero training erroris overfit to the training data and will typically generalize poorly'' (Hastie, Tibshirani, Friedman)



...is still an open problem!

• Increasing the class of allowed functions, it's more likely to find a smooth function with lower norm (complexity): Occam's razor



Mapping Improves Understanding



Representations Make Tasks Easier





Learn Representations



Number of parameters

• Empirical studies: increasing number of parameters doesn't help beyond a certain point



Depth

• Empirical studies: increasing depth tends to always result in some improvement



Regularization: weight decay

- Another way of regularizing is via weight decay (L^2 regularization)
 - Tradeoff between good fitting (small MSE) and small norm (smaller slope, or fewer features with large weights)
 - In statistics, "ridge regression", "Tikhonov regularization"

$$J(\mathbf{w}) = MSE_{train} + \lambda \mathbf{w}^T \mathbf{w}$$



Early stopping...

• Train until the validation set loss starts increasing, and pick the model corresponding to the minimum validation loss



... is a form of regularisation

• Early stopping limits the reachable phase space, and is therefore analogous to L2 regularization (weight decay)

$$J(\mathbf{w}) = MSE_{train} + \lambda \mathbf{w}^T \mathbf{w}$$



Many ways of inserting our biases

- Regularization corresponds to inserting our bias into the algorithm
 - "I know that the solution should not wiggle", "I know that the curvature must not be too large"
- Two models A and B performing the same classification task

$$\circ ~~ \hat{y}^{(A)} = f(\mathbf{w}^{(A)}, \mathbf{x})$$

- $\circ \; \hat{y}^{(B)} = f(\mathbf{w}^{(B)}, \mathbf{x})$
- If the inputs distributions are somehow different, but we know (or want that) the output are related, we can assume that the weights should be similar

$$J(\mathbf{w}) = MSE_{train} + \lambda ||\mathbf{w}^{(A)} - \mathbf{w}^{(B)}||_2^2$$

Parameter sharing

- Simply require parameters are equal
 - If they are equal, you can store only one number in memory (sometimes dramatic memory footprint reduction)



Convolution: a form of averaging

$$s(t)=\int x(a)w(t-a)da$$

- When discretized, integral becomes a sum
 - $\circ x$ input
 - $\circ w$ kernel: specifies how far does the averaging goes
 - *s* feature map

Convolution: a form of averaging

 $S(i,j) = (K*I)(i,j) = \sum \sum I(i-m,j-n)K(m,n)$ m n



Receptive field



Receptive field: deeper = larger



Parameter sharing



Convolutional network

- Convolution \rightarrow nonlinear activation \rightarrow pooling
- Pooling: replace output at a location with a summary statistic
 - e.g., max pooling = report the maximum output in a neighbourhood
 - Helps with invariance for translations



Convolutional networks



Intermediate representations



Morphology of galaxies



Representations of galaxies...



...work pretty well



Semantic representations





Image Recognition

Semantic Segmentation



Object Detection



Instance Segmentation

What about time?

- Convolutional network: process grid of values (e.g. images)
- Recurrent networks: process a sequence of values indicised by a "time" component
 - Language is a sequence
- Parameter sharing crucial to generalize:
 - lengths unseen in training
 - different positions in the sentences
- Without parameter sharing, a network would have to learn all the language rules at each step of the sequence
 - Very impractical
- Both scale very well (thanks to parameter sharing)

Convolutional networks for sequences?

- Could "link" the steps of the sequence via the convolution
- Use the same kernel at each time step
- Shallow: it links only neighbouring time steps

Recurrent network

- Use the same parameter at the same step, $s^{(t)} = f(s^{(t-1)}, heta)$
 - Very deep structure



Unfold the graph

$$s^{(3)}=f(s^{(2)}, heta)=f(f(s^{(1)}, heta), heta)$$



Vast zoology

= An output at each time step, recurrent connections between hidden units



Vast zoology

• An output at each time step, recurrent connections only from the output at one time step to the hidden units at the next time step



Vast zoology

• Recurrent connections between hidden units, that read an entire sequence and then produce a single output



Sequences of images



Real-time segmentation



Graphs Represent Structure



Graph networks

- Represent data as point clouds
- Connect data points with weightdependent connections
- Train the network to find which weights are strongest
 - Learng the connectivity structure of the data



CMS High-granularity calorimeter

- 6 million cells with $\sim 3mm$ spatial resolution, over $600m^2$ of sensors
- Non-projective geometry

Learning representations of irregular particle-detector geometry with distance-weighted graph networks











(d)

Graphs for water simulation


Plug the Physics into the Al



- Data augmentation Li, Dobriban '20
- Loss function penalties concerved quantities
- Architectural design
 - · Approximate symmetries (CNN)

 - Exact symmetries (message passing) (ohen, Welling '19 Ravanbakhsh Weight sharing (group convolutions) Rose 70 '21, '20. Weiler '21 Rose 70 '21, '20. Weiler '21

Enlarge your Dataset

Credit: Bharath Raj

- Parameterization of symmetry preserving functions kondor '18 Maron '18
 Symmetries as constraints Finzi et al '21 Cohen '18
- Irreducible representations Kondor, Thomas '18, Euclus'20 Skerable CNINS Cohen '17, Welling...

Plug the physics into the AI: constraints

$$\hat{y} = f(\mathbf{x}, heta)$$

• Encode physics knowledge (e.g. inconsistency of models) inside the loss function as a penalty term

$$\mathbf{J}(\mathbf{w}) = Loss(y, \hat{y}) + \lambda ||\mathbf{w}||_2^2 + \gamma \Omega(\hat{y}, \Phi)$$



Plug the physics into the AI: network structure

- Equivariance under group transformation can e.g. enforced by convolutional layers
- Some implementations available in pytorch



Plug the Physics into the Al

• Physics-aware differential equations solving



Plug the Physics into the Al

• Several ODE problems now solvable via neural networks



Autoencoders

- Learn the data itself passing by a lower-dimensional intermediate representations
 - Capture data generation features into a lower-dimensional space
- Can use for anomaly detection
- Can sample from the latent space to obtain random samples (generative AI)



Invertible networks

. . .

Inference network



.

. . .

Generative model

Solve inverse problems ("unfolding")

• Correct detector observation noise to recover source distribution



Figure 5: Neural Empirical Bayes for detector correction in collider physics. (a) The source distribution $p(\mathbf{x})$ is shown in blue against the estimated source distribution $q_{\theta}(\mathbf{x})$ in black. (b) Posterior distribution obtained with rejection sampling, with generating source sample \mathbf{x} indicated in red. (c) Calibration curves for each jet property obtained with rejection sampling on 10000 observations. In (a) and (b), contours represent the 68-95-99.7% levels.

Interpretability



Encode sequences

- One-hot encoding for unordered sequences
 - Works e.g. for text

Feature (Color)		One Hot Encoded Vector	Red	Green	Yellow
Red		[1,00]	1	0	0
Green		[0,1,0]	0	1	0
Yellow	One Hot Encoding	[0,0,1]	0	0	1
Green		[0,1,0]	0	1	0
Red		[1,00]	1	0	0

scaler Topics

Encode sequences

- "Yellow" [0,0,1] can be predicted as "0 for red and 0 for green"
 - One-hot-encoded features highly correlated ("multicollinearity")
- Dummy variable trap
 - Drop one of the "dimensions"

Feature (Color)		Red	Green
Red		1	0
Green		0	1
Yellow	One Hot Encoding	0	0
Green Red		0	1
		1	0

Yellow Column dropped to avoid the Dummy Variable Trap



- Capture dependencies and relationships within inputs
 - Mostly in natural language processing and computer vision
- *N* inputs, *N* outputs
 - Allow inputs to interact with eacho other and find out which ones to pay attention to
 - Output is an aggregate of interactions and attention scores
- Useful for:
 - Long-range dependencies: understand complex patterns and dependencies
 - Contextual understanding: assign appropriate weights to important elements in the sequence
 - $\circ~$ Parallel computation: can be computed in parallel \rightarrow efficient and scalable for large datasets.

- Inputs (green) must be represented as: key (orange), query (red), value (purple)
 - Initially, by random reweighting of inputs themselves



1 0

input #2

2 0 2

0

input #3



Self-attention

- Calculate attention score
 - Multiply (dot product) each query with all keys
 - \circ For each query: N keys ightarrow N attention scores

Self-attention



input #2

0

input #3





2 0 2



Pietro Vischia - Statistics for HEP (13th Course on Physics of the LHC, Lisboa, Portugal) - 2024.03.20-22 --- 85 / 108

• Activation function (softmax) of attention scores

Self-attention			
	input #1 1 0 1 0	input #2 0 2 0 2	input #3

- Calculate alignment vectors (yellow), i.e. weighted values
 - Multiply each attention score (blue) by its value (purple)
- Sum alignment vectors to get input for output 1, repeat for 2 and 3

Self-attention





2



Pietro Vischia - Statistics for HEP (13th Course on Physics of the LHC, Lisboa, Portugal) - 2024.03.20-22 --- 87 / 108

Transformers

• The engine behind GPT3



Foundation Models



Translate Problems into Solutions

• Symbolic integration: find the analytic formula for the area of the curve



	Integration (BWD)	ODE (order 1)	ODE (order 2)
Mathematica (30s)	84.0	77.2	61.6
Matlab	65.2	-	-
Maple	67.4	-	-
Beam size 1	98.4	81.2	40.8
Beam size 10	99.6	94.0	73.2
Beam size 50	99.6	97.0	81.0

Reinforcement learning



From Videogames...



...to Physics

• Reward models consistent with the observed quark properties

charges	$\mathcal{Q} = \left(egin{array}{c ccccccccccccccccccccccccccccccccccc$
$\mathcal{O}(1)$ coeff.	$ (a_{ij}) \simeq \begin{pmatrix} -1.975 & 1.284 & -1.219 \\ 1.875 & -1.802 & -0.639 \\ 0.592 & 1.772 & 0.982 \end{pmatrix} (b_{ij}) \simeq \begin{pmatrix} -1.349 & 1.042 & 1.200 \\ 1.632 & 0.830 & -1.758 \\ -1.259 & -1.085 & 1.949 \end{pmatrix} $
VEV, Value	$v_1\simeq 0.224 \;, \qquad \mathcal{V}(\mathcal{Q})\simeq -0.598$
charges	$\mathcal{Q} = \begin{pmatrix} Q_1 & Q_2 & Q_3 & u_1 & u_2 & u_3 & d_1 & d_2 & d_3 & H & \phi \\ \hline 1 & 2 & 0 & -1 & -3 & 1 & -3 & -5 & -4 & 1 & 1 \end{pmatrix}$
$\mathcal{O}(1)$ coeff.	$ (a_{ij}) \simeq \begin{pmatrix} -0.601 & 1.996 & 0.537 \\ -0.976 & -1.498 & -1.156 \\ 1.513 & 1.565 & 0.982 \end{pmatrix} (b_{ij}) \simeq \begin{pmatrix} 0.740 & -1.581 & -1.664 \\ -1.199 & -1.383 & 0.542 \\ 0.968 & 0.679 & -1.153 \end{pmatrix} $
VEV, value	$v_1\simeq 0.158\ ,\qquad {\cal V}({\cal Q})\simeq -0.621$

From perceptron-based networks...

• Matrix multiplication



...to spiking neural networks

- Event-driven computations
 - "when a spike occurs, compute something"



The energy advantage

- Perceptron-based networks: matrix multiplication
 - Sparsity doesn't affect much the throughput and energy consumption
- Spiking neural networks: event-driven computations
 - Sparser inputs require less computations, therefore less time and energy



Encode information with Qubits

- Random bit (Bernoulli random variable) whose description is not governed by classical probability theory but by quantum mechanics
- Not only "because it can take real values in [0,1]": complex numbers as coefficients α and β create interference
 - Interference is not reproducible with classical bits



Represent neural networks

• Qubit operations can represent rather naturally neural networks



 $|\mathrm{in}\rangle\langle\mathrm{in}| = |q_1, q_2, q_3\rangle\langle q_1, q_2, q_3|.$

• Gradient descent exploits intrinsic analytic differentiability of quantum circuits

$$\begin{aligned} \partial_{\mu} \langle \psi(x,\theta) | \sigma_{z} | \psi(x,\theta) \rangle &= \langle 0 | \dots \partial_{\mu} e^{-i\mu\sigma} \dots \sigma_{z} \dots e^{i\mu\sigma} \dots | 0 \rangle \\ &+ \langle 0 | \dots e^{-i\mu\sigma} \dots \sigma_{z} \dots \partial_{\mu} e^{i\mu\sigma} \dots | 0 \rangle \\ &= \langle 0 | \dots (-i\sigma) e^{-i\mu\sigma} \dots \sigma_{z} \dots e^{i\mu\sigma} \dots | 0 \rangle \\ &+ \langle 0 | \dots e^{-i\mu\sigma} \dots \sigma_{z} \dots (i\sigma) e^{i\mu\sigma} \dots | 0 \rangle \\ &= \langle 0 | \dots (1-i\sigma) e^{-i\mu\sigma} \dots \sigma_{z} \dots (1+i\sigma) e^{i\mu\sigma} \dots | 0 \rangle \\ &+ \langle 0 | \dots (1+i\sigma) e^{-i\mu\sigma} \dots \sigma_{z} \dots (1-i\sigma) e^{i\mu\sigma} \dots | 0 \rangle \end{aligned}$$

Pietro Vischia - Statistics for HEP (13th Course on Physics of the LHC, Lisboa, Portugal) - 2024.03.20-22 --- 98 / 108

Need for new paradigma

• If you are interested in Neuromorphic computing or Quantum computing, drop me a line!



Technology readiness?

The Detectors of the Present



The Detectors of the Future

- Optimize/learn by finding the minimum of a function $\mathcal{L}:\mathbb{R}^n
 ightarrow\mathbb{R}$
- Output represents performance on a physics goal or a constraint (e.g. cost)



The Detectors of the Future

• Joint optimization yields in general different solution than optimization of individual features



The Detectors of the Future



Pietro Vischia - Statistics for HEP (13th Course on Physics of the LHC, Lisboa, Portugal) - 2024.03.20-22 --- 103 / 108

MODE

• Creation (11.2020) and rapid expansion of the MODE Collaboration

https://mode-collaboration.github.io/

• Joint effort of particle physicists, nuclear physicists, astrophysicists, and computer scientists

At INFN and Universitá of Padova Dr. Tommaso Dorigo, Dr. Pablo De Castro Manzano, Dr. Lukas Layer, Dr. Giles Strong, Dr. Mia Tosi, and Dr. Hevjin Yarar

At Université catholique de Louvain Dr. Andrea Giammanco, Prof. Christophe Delaere, Mr. Maxime Lagrange, and Dr. Pietro Vischia

At Université Clermont Auvergne, Prof. Julien Donini, and Mr. Federico Nardi (joint with Universitá di Padova)

At the Higher School of Economics of Moscow, Prof. Andrey Ustyuzhanin, Dr. Alexey Boldyrev, Dr. Denis Derkach, and Dr. Fedor Ratnikov At the Instituto de Física de Cantabria, Dr. Pablo Martínez Ruíz del Árbol

At CERN, Dr. Jan Kieseler

At University of Oxford Dr. Atilim Gunes Baydin

At New York University Prof. Kyle Cranmer At Université de Liège Prof. Gilles Louppe

At GSI Dr. Anastasios Belias

At Rutgers University Dr. Claudius Krause

At Uppsala Universitet Prof. Christian Glaser

At TU-München, Prof. Lukas Heinrich and Mr. Max Lamparth

At Durham University Dr. Patrick Stowell

At Lebanese University Prof. Haitham Zaraket

28 people 15 institutions

2022

2019 10 people 4 institutions

for a synergic research plan of potential interest of the JENAS group

T. Dorigo, D. Boumediene, C. Delaere, D. Derkach, J. Donini,

A. Giammanco, R. Rossin, M. Tosi, A. Ustyuzhanin, P. Vischia,

Coordinator
Steering Board

December 3, 2019

Artificial Brain



Artificial General Intelligence?



Pietro Vischia - Statistics for HEP (13th Course on Physics of the LHC, Lisboa, Portugal) - 2024.03.20-22 --- 106 / 108

Not yet.

Pietro Vischia - Statistics for HEP (13th Course on Physics of the LHC, Lisboa, Portugal) - 2024.03.20-22 --- 107 / 108
My efforts in the next years

- Exploration of high-dimensional spaces via gradient descent, eventually powered by quantum algorithms
- Realistic neurons (spiking networks on neuromorphic circuits)
- Applications to experiment design and to heart diseases
- If you are interested in collaborating, drop me an email!