Particle Identification with Autoencoder Neural Networks

Daniel Lersch & Sean Dobbs

June 19, 2020

Daniel Lersch (FSU)

JLUO AI Meeting

June 19, 2020 1 / 24

Overview

- 1. Very short introduction to machine learning
- 2. Definition and properties of (deep) neural networks
- 3. Application of autoencoders for particle identification at GlueX











- Any machine learning algorithms "learn" patterns / actions from a given data set by setting its internal parameters appropriately
- The parameter adjustment is done during training



- Any machine learning algorithms "learn" patterns / actions from a given data set by setting its internal parameters appropriately
- The parameter adjustment is done during training



- Any machine learning algorithms "learn" patterns / actions from a given data set by setting its internal parameters appropriately
- The parameter adjustment is done during training



- Any machine learning algorithms "learn" patterns / actions from a given data set by setting its internal parameters appropriately
- The parameter adjustment is done during training



clustering algorithm taken from: wikinedia

By Chire - Own work, CC BY-SA 3.0, https://commons.wikimedia.org/w/index.php?curid=1708708

- Any machine learning algorithms "learn" patterns / actions from a given data set by setting its internal parameters appropriately
- The parameter adjustment is done during training



- Any machine learning algorithms "learn" patterns / actions from a given data set by setting its internal parameters appropriately
- The parameter adjustment is done during training
- NOTE: The algorithms behavior / performance highly depends on the provided training data



- Any machine learning algorithms "learn" patterns / actions from a given data set by setting its internal parameters appropriately
- The parameter adjustment is done during training
- NOTE: The algorithms behavior / performance highly depends on the provided training data
- Your algorithm is useless if the training data significantly differs from the data you are trying to analyze (However, there is some room for variations...)



The Multilayer Perceptron



- Most popular example for machine learning algorithms
- Belongs to the class of feedforward neural networks
- Architecture: Hidden layers with a set of neurons

Daniel Lersch (FSU)

The Multilayer Perceptron



- Most popular example for machine learning algorithms
- Belongs to the class of feedforward neural networks
- Architecture: Hidden layers with a set of neurons

Daniel Lersch (FSU)

A single Neuron



• Basic ingredients: Information from previous neurons, weights, bias and activation function

A single Neuron



- Basic ingredients: Information from previous neurons, weights, bias and activation function
- Activation function plots taken from Mustafa Mustafas lecture at the: deep learning for science school 2019

Daniel Lersch (FSU)

A single Neuron



- Basic ingredients: Information from previous neurons, weights, bias and activation function
- Activation function plots taken from Mustafa Mustafas lecture at the: deep learning for science school 2019

Daniel Lersch (FSU)

The Universal Approximation Theorem for Neural Networks

"a single hidden layer neural network with a linear output unit can approximate any continuous function arbitrarily well, given enough hidden units" -- Hornik, 1991, http://zmiones.com/static/statistical-learning/hornik-nn-1991.pdf

This, of course, does not imply that we have an optimization algorithm that can find such a function. The layer could also be too large to be practical.



Fig. credit towardsdatascience.com/can-neural-networks-really-learn-any-function-65e106617fc6

Screenshot taken from Mustafa Mustafas lecture at the: deep learning for science school 2019

 \Rightarrow Similarly formulated in 1990 by the **Stone-Weierstrass-Theorem**

"[...] there are no nemesis functions that cannot be modeled by neural networks"

Training a (deep) Neural Network

Bottom left picture taken from here

Bottom right picture taken from Mustafa Mustafas talk at the: deep learning for science school 2019



Gradient Descent:
$$w_{k+1} = w_k - \eta \nabla L(w_k) = w_k - \frac{\eta}{m} \sum_{i=1}^m \nabla L(x_i, w_k)$$
 (1)

- It is all about updating the weights w_k
 - L: Loss function \propto prediction truth
 - η : Learning rate (Common choice: 1/N(Training epochs))
 - m: Batch size
- Problem: $\nabla L = 0$ (vanishing gradient)
- Different variants of stochastic gradient descent (L-BFGS, Adam, SGD,...)

Daniel Lersch (FSU)

 Want to enable network to abstract / generalize on unknown data AND avoid overfitting (i.e. avoid that network reproduces features from training data only)



Picture taken from Brenda Ngs introductory talk at the: deep learning for science school 2019

- Want to enable network to abstract / generalize on unknown data AND avoid overfitting (i.e. avoid that network reproduces features from training data only)
- Validation Data: Part of training data that is NOT used to update internal parameters¹, but used to determine when training is complete

¹This data is "unseen" by the algorithm during the training stage

- Want to enable network to abstract / generalize on unknown data AND avoid overfitting (i.e. avoid that network reproduces features from training data only)
- Validation Data: Part of training data that is NOT used to update internal parameters¹, but used to determine when training is complete



Picture taken from Mustafa Mustafas talk at the: deep learning for science school 2019

¹This data is "unseen" by the algorithm during the training stage

Daniel Lersch (FSU)

- Want to enable network to abstract / generalize on unknown data AND avoid overfitting (i.e. avoid that network reproduces features from training data only)
- Validation Data: Part of training data that is NOT used to update internal parameters¹, but used to determine when training is complete



Picture taken from Mustafa Mustafas talk at the: deep learning for science school 2019

¹This data is "unseen" by the algorithm during the training stage

Daniel Lersch (FSU)

• "Classical" Machine learning

- "Classical" Machine learning
 - ► Features (inputs for your algorithm) are obtained after some pre-processing (calibration, analysis cuts, variable selection,...) → Also known as feature-engineering
 - By pre-processing, you already encode information into the data
 - Moderate model size (e.g 1-2 hidden layers in a neural network)



- "Classical" Machine learning
 - ► Features (inputs for your algorithm) are obtained after some pre-processing (calibration, analysis cuts, variable selection,...) → Also known as feature-engineering
 - By pre-processing, you already encode information into the data
 - Moderate model size (e.g 1-2 hidden layers in a neural network)
- Deep learning
 - Still machine learning, but uses neural networks only
 - Leave out (certain) pre-processing steps \rightarrow Let the model do the work for you
 - \blacktriangleright The neural network becomes deep \rightarrow Pre-processing is basically done in extra hidden layers



Picture taken from here

Daniel Lersch (FSU)

• "Classical" Machine learning

- ► Features (inputs for your algorithm) are obtained after some pre-processing (calibration, analysis cuts, variable selection,...) → Also known as feature-engineering
- By pre-processing, you already encode information into the data
- Moderate model size (e.g 1-2 hidden layers in a neural network)
- Deep learning
 - Still machine learning, but uses neural networks only
 - Leave out (certain) pre-processing steps \rightarrow Let the model do the work for you
 - \blacktriangleright The neural network becomes deep \rightarrow Pre-processing is basically done in extra hidden layers
- Deep learning is not trivial, but fortunately there are many frameworks
 - Pytorch
 - Keras
 - Tensorflow (used for the work presented here)

► ...



Dense layer with n neurons (n scales with height)

- Consists of two parts: encoder + decoder \rightarrow Symmetric sandwich architecture
- Train autoencoder for: Data In \approx Data Out
- But: Number of neurons needs to decrease in encoding part and increase in decoding part → Otherwise no effect ↔ Identity function

Daniel Lersch (FSU)



Dense layer with n neurons (n scales with height)

- Consists of two parts: encoder + decoder \rightarrow Symmetric sandwich architecture
- Train autoencoder for: Data In \approx Data Out
- But: Number of neurons needs to decrease in encoding part and increase in decoding part → Otherwise no effect ↔ Identity function

Daniel Lersch (FSU)



- Consists of two parts: encoder + decoder \rightarrow Symmetric sandwich architecture
- Train autoencoder for: Data In \approx Data Out
- But: Number of neurons needs to decrease in encoding part and increase in decoding part → Otherwise no effect ↔ Identity function

Daniel Lersch (FSU)



Dense layer with n neurons (n scales with height)

- ullet Consists of two parts: encoder + decoder ightarrow Symmetric sandwich architecture
- Train autoencoder for: Data In \approx Data Out
- But: Number of neurons needs to decrease in encoding part and increase in decoding part → Otherwise no effect ↔ Identity function

Daniel Lersch (FSU)

Autoencoders: Compression and Decompression



Features:

- momentum
- theta
- ddEdx(CDC)
- dE(BCAL 0)
- dE(BCAL all)
- ddEdx(FDC)
- dE(FCAL)
- E9E25 (FCAL)
- E1E9 (FCAL)

Autoencoders: Compression and Decompression



Compression / dimensional reduction from $9 \rightarrow 4$ dimensions

Autoencoders: Compression and Decompression



Daniel Lersch (FSU)

Feature correlations similar to training data



Daniel Lersch (FSU)





Daniel Lersch (FSU)

Feature correlations similar to training data



Daniel Lersch (FSU)

Feature correlations similar to training data



Daniel Lersch (FSU)

Autoencoders: (Current) Realization at GlueX



 Train autoencoder on simulated single particle tracks → One autoencoder per particle species and per charge (as apposed to "classic" machine learning)

Autoencoders: (Current) Realization at GlueX



- Train autoencoder on simulated single particle tracks → One autoencoder per particle species and per charge (as apposed to "classic" machine learning)
- Autoencoder outputs help to understand model performance

Autoencoders: (Current) Realization at GlueX



- Train autoencoder on simulated single particle tracks → One autoencoder per particle species and per charge (as apposed to "classic" machine learning)
- Autoencoder outputs help to understand model performance
- Anomaly detector (aka classifer) allows to translate the decoded data into a likelihood

Reference Models and Analysis Procedures

- It is very helpful to have a reference model², running in parallel to your machine / deep learning analysis
 - \Rightarrow Helps to understand / debug the model your are developing

²Preferably not machine learning based

Reference Models and Analysis Procedures

- It is very helpful to have a reference model², running in parallel to your machine / deep learning analysis
 - \Rightarrow Helps to understand / debug the model your are developing
- In this Analysis use:
 - 1. E/p-cut: Uses energy deposit E in either BCAL / FCAL and the momentum $p \rightarrow$ Not machine learning based and straight forward to apply
 - 2. "Classic" machine learning approach: Neural network trained on leptons vs. pions (one for each charge)

²Preferably not machine learning based

Reference Models and Analysis Procedures

- It is very helpful to have a reference model², running in parallel to your machine / deep learning analysis
 - \Rightarrow Helps to understand / debug the model your are developing
- In this Analysis use:
 - 1. E/p-cut: Uses energy deposit E in either BCAL / FCAL and the momentum $p \rightarrow$ Not machine learning based and straight forward to apply
 - 2. "Classic" machine learning approach: Neural network trained on leptons vs. pions (one for each charge)
- Use Bayesian formalism to combine predictions from pion- and leptonautoencoders into a likelihood

²Preferably not machine learning based

Lepton Identification Performance on simulated single e^+ - $/\pi^+$ -Tracks



Daniel Lersch (FSU)

Lepton Identification Performance on simulated single e^{-}/π^{-} -Tracks



Daniel Lersch (FSU)

Application on GlueX Lepton Data



- Analyzed small sample of GlueX e⁺e⁻-trees
- Did not tune cuts
- Found: $\Phi \rightarrow e^+e^-$ and $J/\psi \rightarrow e^+e^-$

- Analyzed small sample of GlueX e⁺e⁻γ-trees
- Used rather tight cuts: Networks \gtrsim 70%, $E/p \in [0.8, 1.1]$
- Found: $\pi^0 \to e^+ e^- \gamma$ and $\eta \to e^+ e^- \gamma$

Application on GlueX Lepton Data



- Analyzed small sample of GlueX e^+e^- -trees
- Did not tune cuts
- Found: $\Phi \rightarrow e^+e^-$ and $J/\psi \rightarrow e^+e^-$

- Analyzed small sample of GlueX e⁺e⁻γ-trees
- Used rather tight cuts: Networks \gtrsim 90%, $E/p \in [0.85, 1.05]$
- Found: $\pi^{0} \rightarrow e^{+}e^{-}\gamma$ and $\eta \rightarrow e^{+}e^{-}\gamma$

Predicted Lepton Features



• Remeber: $E/p \sim 1$ for leptons and predominantly ~ 0 for pions

Autoencoder reconstructs E/p by taking correlations of provided features into account

- \Rightarrow basically a multidimensional fit
- \Rightarrow similar to a kinematic fit, but on a single track level!
- Cuts on predicted E/p allow for background suppression

• Very important - not only for you, but also for the analyzer who is using your model

- Very important not only for you, but also for the analyzer who is using your model
- Autoencoder allows to formulate residuals: Network Input Network Output

- Very important not only for you, but also for the analyzer who is using your model
- Autoencoder allows to formulate residuals: Network Input Network Output
- Shown below: $\chi^2(\text{Residual}) = \sum_{i=1}^{N} \left(\text{Network Input}[i] \text{Network Output}[i] \right)^2$ and a few residuals BEFORE a bug-fix



• χ^2 and residuals expected to be pprox 0

- Very important not only for you, but also for the analyzer who is using your model
- Autoencoder allows to formulate residuals: Network Input Network Output
- Shown below: $\chi^2(\text{Residual}) = \sum_{i=1}^{N} \left(\text{Network Input}[i] \text{Network Output}[i] \right)^2$ and a few residuals BEFORE a bug-fix



• Turned out that energy deposits of BCAL layers >=2 are -1 in DSelector

- Very important not only for you, but also for the analyzer who is using your model
- Autoencoder allows to formulate residuals: Network Input Network Output
- Shown below: $\chi^2(\text{Residual}) = \sum_{i=1}^{N} \left(\text{Network Input}[i] \text{Network Output}[i] \right)^2$ and a few residuals AFTER a bug-fix



Current model uses BCAL pre-shower and sum over all Layers only

• χ^2 and residuals are pprox 0 ightarrow Deviations due to pions and poorly reconstructed events

- Very important not only for you, but also for the analyzer who is using your model
- Autoencoder allows to formulate residuals: Network Input Network Output
- Shown below: χ^2 (Residual) = $\sum_{i=1}^{N} \left(\text{Network Input}[i] \text{Network Output}[i] \right)^2$ and a few residuals AFTER a bug-fix



• Current model uses BCAL pre-shower and sum over all Layers only

• χ^2 and residuals are pprox 0 ightarrow Deviations due to pions and poorly reconstructed events

Inspecting the Decoder (first generation autoencoder)



Inspecting the Decoder (first generation autoencoder)



Inspecting the Decoder (first generation autoencoder)



binned Histogram

Generating Lepton (like) Data (first generation autoencoder)

Left: Original Data / Right: Generated Data



Daniel Lersch (FSU)

Moving forward: Constraint hypothesis fitting with Autoencoders



- Analogue to kinematic fitter: Include constraints related to particle hypothesis
- For leptons: Energy deposits in calorimeter $E(cal) \sim$ Momentum p
- Right column shows network predictions, using the lepton constraints ⇒ Improvement in resolution!

Daniel Lersch (FSU)

Summary and Outlook

- Used autoencoder neural networks for PID in GlueX
- Obtained promising results
- Might be a bit too much for lepton ID, but:
 - the goal is to design a generic algorithm
 - want to proof and develop the concept on an "easy" case
- Besides simple classification, autoencoder networks...
 - ... support the user regarding debugging and understanding the model
 - ... interpret the data with respect to a given particle hypothesis (analogue to a kinematic fitter)
 - ... are easy to train due to the constraint architecture (used tied weights here)
 - ... de-noise the data (not shown here today)
 - \dots are somewhat more robust with respect to training data / real data mismatch
- $\bullet\,$ The price one has to pay is the model size: $\sim 10^3~\text{vs} \sim 10^2$ parameters