

Boyang Yu, Nikolai Hartmann, Luca Schinnerl, Thomas Kuhr Joint Seminar of Particle Physics Groups, July 13th, 2022







Contents

Contents

- Motivation
- Machine Learning Basics
 - Categories
 - Linear Regression
 - Gradient Descent
 - Logistic Regression
 - Neural Network
 - Activation Functions
- Graph Attention Networks
 - Data Structure
 - GAT and GAP
 - Network Structure
- Weighting Methods
- Performances and Generalization
- Summary





Motivation

Current Monte Carlo Simulation data flow







Motivation

Current Monte Carlo Simulation data flow





Motivation

Current Monte Carlo Simulation data flow



Previous Works:

- **PhD Thesis**: Hadronic Tag Sensitivity Study of $B \to K^{(*)}v\bar{v}$ and Selective Background Monte Carlo Simulation at Belle II, James Kahn, 2019
- **Talk**: Selective background Monte Carlo simulation at Belle II, James Kahn, CHEP 2019





- Easier to achieve
- Beyond human limits of knowledge and observation
- Simulation with high efficiency
- Large datasets available











Our Task:

- **In:** Particle decay generated by Monte Carlo •
- **Out:** Whether it will pass the skim afterwards ٠

Supervised, Classification





Machine Learning Basics

Start with Linear Regression



Estimation:

$$\hat{y} \stackrel{\text{\tiny def}}{=} b + w_1 x + w_2 x^2 + w_3 x^3 \stackrel{\text{\tiny def}}{=} \overrightarrow{W} \cdot \overrightarrow{X} + b = W^i X_i + b$$
with $\vec{X} = \begin{pmatrix} x \\ x^2 \\ x^3 \end{pmatrix}$



Machine Learning Basics

Start with Linear Regression



Estimation:

$$\hat{y} \stackrel{\text{\tiny def}}{=} b + w_1 x + w_2 x^2 + w_3 x^3 \stackrel{\text{\tiny def}}{=} \overrightarrow{W} \cdot \overrightarrow{X} + b = W^i X_i + b$$
with $\vec{X} = \begin{pmatrix} x \\ x^2 \\ x^3 \end{pmatrix}$

Loss: the difference between estimation and truth

$$L \stackrel{\text{\tiny def}}{=} (\hat{y} - y)^2$$

Cost: the average of losses over all the samples

$$J(W,b) \stackrel{\text{\tiny def}}{=} \frac{1}{2m} \sum_{i=1}^{m} L_i = \frac{1}{2m} \sum_{i=1}^{m} (\hat{y}(x_i) - y(x_i))^2$$





Machine Learning Basics

Start with Linear Regression



Estimation:

$$\hat{y} \stackrel{\text{def}}{=} b + w_1 x + w_2 x^2 + w_3 x^3 \stackrel{\text{def}}{=} \overrightarrow{W} \cdot \overrightarrow{X} + b = W^i X_i + b$$
with $\vec{X} = \begin{pmatrix} x \\ x^2 \\ x^3 \end{pmatrix}$

Loss: the difference between estimation and truth

$$L \stackrel{\text{\tiny def}}{=} (\hat{y} - y)^2$$

Cost: the average of losses over all the samples

$$J(W,b) \stackrel{\text{\tiny def}}{=} \frac{1}{2m} \sum_{i=1}^{m} L_i = \frac{1}{2m} \sum_{i=1}^{m} (\hat{y}(x_i) - y(x_i))^2$$

Training: find proper parameters **W**, *b* to minimize the cost -> Best model





Machine Learning Basics

Start with Linear Regression



Estimation:

$$\hat{y} \stackrel{\text{\tiny def}}{=} b + w_1 x + w_2 x^2 + w_3 x^3 \stackrel{\text{\tiny def}}{=} \overrightarrow{W} \cdot \overrightarrow{X} + b = W^i X_i + b$$
with $\vec{X} = \begin{pmatrix} x \\ x^2 \\ x^3 \end{pmatrix}$

Loss: the difference between estimation and truth

$$L \stackrel{\text{\tiny def}}{=} (\hat{y} - y)^2$$

Cost: the average of losses over all the samples

$$J(W,b) \stackrel{\text{\tiny def}}{=} \frac{1}{2m} \sum_{i=1}^{m} L_i = \frac{1}{2m} \sum_{i=1}^{m} (\hat{y}(x_i) - y(x_i))^2$$

Training: find proper parameters W, b to minimize the cost -> Best model Complexity: $O(n^2)$





Training: Gradient Descent

Update parameters in the opposite direction of the gradient of loss function in each step







Training: Gradient Descent

Update parameters in the opposite direction of the gradient of loss function in each step





Various of algorithms based on GD:

- Batching: MiniBatchGD, StochasticGD, ...
- Momentum: Nesterov, Adam, …



Classification: Logistic Regression









Classification: Logistic Regression

Probability of passing exam versus hours of studying



• Regression with boundary [0,1]





Machine Learning Basics

Classification: Logistic Regression

Probability of passing exam versus hours of studying



- Regression with boundary [0,1]
- Represent the probability of an event occurring



Machine Learning Basics

Classification: Logistic Regression

Probability of passing exam versus hours of studying



- Regression with boundary [0,1]
- Represent the probability of an event occurring
- Generated by logistic function:

$$g(z) = \frac{1}{1+e^{-z}}$$
 with $z = \vec{w} \cdot \vec{x} + b$ called "logit"





Machine Learning Basics

Classification: Logistic Regression

Probability of passing exam versus hours of studying



- Regression with boundary [0,1]
- Represent the probability of an event occurring
- Generated by logistic function:

$$g(z) = \frac{1}{1+e^{-z}}$$
 with $z = \vec{w} \cdot \vec{x} + b$ called "logit"



Classification achieved by introducing a threshold



Machine Learning Basics

Neural network









ÿ



8/19



8/19



Activation Functions: bring non-linearities







Activation Functions: bring non-linearities





- Sigmoid functions (Logistic, Tanh)
 - Bounded ideal for classification
 - Gradient vanish easily harmful for hidden layers



Activation Functions: bring non-linearities





- Sigmoid functions (Logistic, Tanh)
 - Bounded ideal for classification
 - Gradient vanish easily harmful for hidden layers
- (Partially) Linear activations
 - Fast to calculate
 - Allowing large values
 - Not differentiable at 0 some solutions available





- **D** Each particle (each **Node**)
 - PDG ids
 - > 8 Features: Production time, Energy, Position (3d), Momentum (3d)
- Labels per event: Pass/Fail after the reconstruction of B decays (FEI skims)
- Other event level **attributions** for further analysis: e.g. M_{bc} , etc.





 $\vec{\alpha}_{13}$

 $ec{h}_5$

 \vec{h}_4

 $ec{h}_6$



- Update the node \vec{h}_1 by adding up the information from all of its neighbours and itself
- Each pair has an attention weight α_{ij} , calculated • from the information of both nodes \vec{h}_i and \vec{h}_i

 $\vec{\mathbf{a}}$

 $\mathbf{W}\vec{h}_i$

 $\mathbf{W}\vec{h}_i$



Graph Attention Networks



LUDWIG-

MÜNCHEN



- Update the node \vec{h}_1 by adding up the • information from all of its neighbours and itself
- Each pair has an attention weight α_{ij} , calculated ٠ from the information of both nodes \vec{h}_i and \vec{h}_i
- Several sets of attention weights are trained • parallelly to extract different aspects of correlation between nodes



Graph Attention Network (GAT) α_{ij} h_2 softmax concat/avg $\vec{\alpha}_{13}$ $ec{h}_6$ ā h_5 $\mathbf{W}\vec{h}_i$ $\mathbf{W}\vec{h}_i$



- Update the node \vec{h}_1 by adding up the information from all of its neighbours and itself
- Each pair has an attention weight α_{ij} , calculated from the information of both nodes \vec{h}_i and \vec{h}_j
- Several sets of attention weights are trained parallelly to extract different aspects of correlation between nodes

Global Attention Pooling (GAP)

- Aggregate the information from all the nodes with learned attention weight for each node
 - -> Deconstruction of graph structure





Final Architecture: GAT …







Final Architecture: GAT …


























Smart Background Simulation with Graph Attention Networks and Weighting Methods

Event Weighting Methods



	Sampling Method	Reweighting Method
Use of NN output	As probability to keep event randomly	As score for selection according to fix threshold
Weight		
Loss to train NN		



Smart Background Simulation with Graph Attention Networks and Weighting Methods

Event Weighting Methods



	Sampling Method	Reweighting Method
Use of NN output	As probability to keep event randomly	As score for selection according to fix threshold
Weight	Inverse of NN output	Decided with the help of another classifier
Loss to train NN		



Smart Background Simulation with Graph Attention Networks and Weighting Methods

Event Weighting Methods



	Sampling Method	Reweighting Method
Use of NN output	As probability to keep event randomly	As score for selection according to fix threshold
Weight	Inverse of NN output	Decided with the help of another classifier
Loss to train NN	Speedup	Binary cross entropy





--Improvement of computation time to produce the same

effective number of events with the help of NN filter:

Speedup:
$$s = \frac{t_{no_filter}}{t_{filter}}$$

Effective Sample Size:
$$N_{eff} = \frac{(\sum \omega_i)^2}{\sum \omega_i^2}$$





Metric: Binary Cross Entropy

-- Amount of information provided by the prediction





Performances





Performances

Maximum speedup

Bias



6.5

Potential residual bias

2.0

No bias

MINU



Performances on another skim with sampling









Results:

- Graph Attention Networks can learn the information from particle decay trees and help to simulate skims with low retention rate at early stage
- Bias is avoided with sampling method while a speedup of factor 2 can still be maintained.
- Reweighting method can reach much higher speedup up to 6.5 but will still have some bias in the variables that are not used in the training of the extra classifier.
- SmartBKG can be invested to other skims



Thank You for your Attention

Boyang Yu, Nikolai Hartmann, Luca Schinnerl, Thomas Kuhr Joint Seminar of Particle Physics Groups, July 13th, 2022





Backup



Smart Background Simulation with Graph Neural Networks



Backup: Network Structures



Smart Background Simulation with Graph Neural Networks

Backup: Network Structures



LUDWIG-MAXIMILIANS-UNIVERSITÄT

MÜNCHEN







Backup: Network Structures

Quantitative Studies



Backup: Evaluations

Comparison

Parameters:

LUDWIG-MAXIMILIANS-UNIVERSITÄT

MÜNCHEN

- n heads = 4٠
- $n_{layers} = 6$ •
- n_units = 128 •
- $batch_size = 128$ ٠
- n train = 0.9M•
- n_val = 0.1M •
- n test = 0.5M •

Loss:

Entropy ٠

EarlyStopping:

- patience = 3٠
- delta = 1e-5





 GCNsep GATsep

GATgen

10

12

1.0

0.8

	GCN(sep)	GAT(sep)	GAT(gen)	GAT+GAP(gen)
TrainingTime	3619.46s	4047.47s	3471.48s	5049.81s
AUCValues	0.90831	0.90937	0.90891	0.91216

Backup: Evaluations

Grid Search

LM

LUDWIG-<u>MAXIMIL</u>IANS-

MÜNCHEN

UNIVERSITÄT



Batch- size	Number of units	AUC	Training Time
<mark>128</mark>	<mark>128</mark>	<mark>0.9117</mark>	<mark>5205</mark>
256	32	0.9105	4061
256	128	0.9105	2666
<mark>512</mark>	<mark>32</mark>	<mark>0.9117</mark>	<mark>3568</mark>
512	128	0.9115	2228
<mark>1024</mark>	<mark>32</mark>	<mark>0.9115</mark>	<mark>1716</mark>
1024	256	0.9102	3556

Network Sizes

Best Combinations

# Units	# Parameters
32	120,527
64	459,951
128	1,808,495
256	7,184,367
512	28,651,247



Sampling

- Previous method: Cut according to neural network outputs
- Problem: Inevitable bias
- Our method: Sampling with probability given by neural network outputs
- Problem: Statistical uncertainty







Motivations

Implementation of NN filters





Event Weighting Methods

Sampling Method:







Event Weighting Methods

Reweighting Method:



12/15

Studied reweighters:

- GBDT Reweighting
- Histogram Reweighting



Event Weighting Methods

Reweighting Method:

- Train a Gradient Boosting Decision Tree (GBDT) classifier with some event level variables to distinguish between True-Positve events and False-Negative events
- GBDT Reweighting: use the outputs of the classifier directly:

$$W = \frac{1}{p_{clf}} = \frac{1}{p_{TP}/p_{TP+FN}} = \frac{p_{pass_skim}}{p_{TP}}$$

 Histogram Reweighting: compare the score histogram of all the events that can pass the skim (True-Positive + False-Negative) with the score histogram of True-Positives to give each bin of score a scaling factor:

$$w = w_{bin_i|p_{clf} \in bin_i} = \frac{H_{pass_skim,i}}{H_{TP,i}}|_{p_{clf} \in bin_i}$$

Skim NN	Positive	Negative
Pass	True- <mark>Positive</mark> (TP)	False- <mark>Negative</mark> (F <mark>N</mark>)
Fail	False-Positive (FP)	True- <mark>Negative</mark> (TN)





Sampling



Relative statistical uncertainty and effective sample size

Variable	Formula	Remark
NN outputs / Probabilities to pass	$\{p_i\}$	'i' refers to each event in the whole sample (batch)
Weights	$\{\omega_i\} = \left\{\frac{1}{p_i}\right\}$	Infinities (at $p_i = 0$) are excluded and set to 0 Avoid the bias by construction
Relative statistical uncertainty	$S = \frac{\sqrt{\sum \omega_i^2 p_i}}{\sum \omega_i p_i}$	$\sum \omega_i^2 p_i = \sum \omega_i$ $\sum \omega_i p_i = N$ Here consider only passed events (label = 1)
Effective sample size	$N_{eff} = \frac{1}{S^2}$	Number of events needed to reach the same statistical uncertainty without sampling



Backup: Speedup

Speedup rate



Variable	Formula	Remark
Skim retention rate	r = 0.05	Probability to pass the skim process
Times of different phases in ms	$t_{gen} = 0.08$ $t_{NN} = 0.63$ $t_{SR} = 97.04$	Taken from previous studies
Effective number of events after sampling	$n_{+} = \sum p_{i}$ $n_{-} = \sum (1 - p_{i})$	$\{p_i\}$ will be devided into two subsets where the events will/won't pass the skim process
Time consuming with NN filter	$t_{+} = [n_{TP}r + n_{FP}(1 - r)](t_{gen} + t_{NN} + t_{SR})$ $t_{-} = [n_{FN}r + n_{TN}(1 - r)](t_{gen} + t_{NN})$	Positive/Negative: Result of sampling True/False: Result of sampling == skim process
Time consuming without NN	$t_0 = N_{eff} (t_{gen} + t_{NN})$	To reach the same statistical uncertainty
(Inverse) Speedup rate	$R = \frac{t_+ + t}{t_0}$	The lower the better



LUDWIG-MAXIMILIANS-UNIVERSITÄT

MÜNCHEN

Performances of different training losses Speedup always as validation loss!

Entropy: Binary cross entropy between NN outputs $\{p_i\}$ and skim results $\{y_i \in \{0,1\}\}$ **Speedup**: (Inverse) Speedup rate *R* **Entropeedup**: Training-epoch-dependent combination of Entropy and Speedup

Entropeedup = $e^{-epoch/10}Entropy + Speedup$







Metric: Speedup

--Improvement of computation time to produce the same effective number of events with the help of NN filter:

Speedup =
$$\frac{t_{no_filter}}{t_{filter}} \sim \text{Effective Sample Size} = \left(\frac{(\sum \omega_i p_i)^2}{\sum \omega_i^2 p_i}\right)$$

Robustness:







