

Introduction into Deep Learning

Andrey Ustyuzhanin Yandex School of Data Analysis, Higher School of Economics



INFN school of statistics, Ischia, May 2017



Intro into Intro

 \searrow

What is Yandex

Public company since 2011

In Q3 2016 search share across different platforms were approximately¹: 64% on desktop 38% on Android 42% on iOS

Source: Liveinternet.ru 2012-December, 2016; includes desktop and mobile

¹ Based on company estimates, as provided on Q3 2016 earnings call Andrey Ustyuzhanin





Our group

Yandex School of Data Analysis (YSDA) - non-commercial educational organisation; Research group at Yandex School of Data Analysis

2 physicists (PhD), 8 data scientists (6 of them are graduate/undergraduate students from MIPT, MSU, HSE)

Laboratory of Methods for Big Data Analysis, HSE NRU YSDA is member of HEP collaborations:

- CERN: LHCb (since 2015), SHiP (since 2014)
- > CRAYFIS (since 2015), OPERA

Andrey Ustyuzhanin

SCHOOL OF DATA ANALYSIS





Working group directions

Research & Development

- Addressing Physics challenges by means of Machine Learning a) experiment design, b) advanced modelling, c) fast simulation, d) tracking Education
- Machine Learning Courses (YSDA, Imperial College London, Helsinki CSC) Summer Schools on Machine Learning (2015-...)

Outreach

- Masterclasses, Data&Science (https://events.yandex.ru/events/ds/), > Hackathones on Data Science & Machine Learning (<u>http://bit.ly/2lxUWCO</u>)



Why going deep?

- > image recognition
- > text recognition
- > natural text translation
- > voice recognition
- > Go
- > unmanned vehicle control
- > ...(you name it)

Andrey Ustyuzhanin

Physics examples

- > Higgs boson exotic decay (ATLAS, simulated dataset)
- > Jet identification (ATLAS, CMS)
- > Muon tracking (CRAYFIS)
- > Data Certification (LHCb, CMS)



Artificial Neural Networks

Artificial Neural Network history **Origin:**

1943, W. McCulloch & W. Pitts — first computational model 1949, D. Hebb — method for learning 1954, W. Clark — first simulations

Further development:

1970, Linnainmaa, *Backpropagation* in the modern form 1979, Fukushima The earliest convolutional networks 1985, Rumelhart, Hinton, and Williams, backpropagation in neural networks could yield interesting representations 1989, LeCun,

convolutional networks + *backpropagation* = LeNet





Alstages

1965 The first deep network architecture, A. Ivakhnenko (USSR)

1980s, Al winter - in applications simple (linear) models models (in particular, decision trees, SVM) clearly demonstrate better performance than Neural Nets

2006 — till now Neural networks Renaissance





Single Neuron

Neuron can be activated or not (h = 1, h = 0) by signals from receptors or other neurons.

in the simplest case:

$$h=\Theta(\sum w_i x_i + b)$$

> $\Theta(x)$ — Heaviside function (0, if x \leq 0, 1 otherwise)

> non-smooth \Rightarrow hard to optimize

can we use other functions?







Activation functions

Sigmoid:

$$f(x) = \frac{1}{1 + e^{-x}}$$

ReLU - rectifier linear unit (biologically inspired)

 $f(x) = \max(0, x)$

Softplus (smoothed version of ReLU)

$$f(x) = \ln(1 + e^x)$$



Artificial Neural Network



hidden layer

How many parameters are there between hidden layer 1 and hidden layer 2? (right side)

Andrey Ustyuzhanin



hidden layer 1 hidden layer 2



Computations in neural network

input:

 $(x_1, x_2, ..., x_d)$

Activations of hidden layer:

$$h_i = f(\sum_{j=1}^{n_{inputs}} W_{ij} x_j + b_i)$$

Activations of output layer:

$$o_i = f(\sum_{j=1}^{n_{hidden}} \widetilde{W}_{ij} h_j + \widetilde{b}_i)$$

Andrey Ustyuzhanin



hidden layer

$$\begin{cases} h_{linear} = Wx + b \\ h = f(h_{linear}) \end{cases}$$
$$o_{linear} = \widetilde{W}h + \widetilde{b} \\ o = f(o_{linear}) \end{cases}$$



Several layers (Multi Layer Perceptron)

in matrix form, *n* hidden layers:

$$x_{1} = f(W_{0}x_{0} + b_{0})$$

$$x_{2} = f(W_{1}x_{1} + b_{1})$$

$$x_{3} = f(W_{2}x_{2} + b_{2})$$

 $\hat{\mathbf{y}} = f(W_n x_n + b_n) - output$

. . .

should compare with *y* - true labels corresponding to *X*.

Andrey Ustyuzhanin



hidden layer 1 hidden layer 2



Training Neural Net

Mathematically, neural network is a function with many parameters, so we use smooth optimization. Output of network can be a single value (in case of one neuron) or several.

The only thing needed is to have some smooth loss function for optimization

$$l(y, \hat{y}) = \sum_{i} (\hat{y}_{i} - y_{i})^{2}$$







Gradient Descent

Problem: find *W* to minimize *l*. Solution: use gradient descent (repeat until convergence.)

$$w \leftarrow w - \eta \frac{\partial l}{\partial w}$$

η is a step size (also called *shrinkage*, *learning rate*)





Other kinds of optimizers

- Stochastic gradient descent:
- calculate error using subsample of training set
- Stochastic gradient descent with momentum:
- moves towards "overall gradient direction", not just current gradient

AdaGrad:

decreases learning rate individually for each parameter in proportion to sum of it's gradients so far

RMSProp:

> makes sure all gradient steps have ~ same magnitude (by keeping moving) average of step magnitude)









Regularizations

Many parameters → overfitting

To avoid too complicated models <u>regu</u> function:

$$L_1 = ||W||_1, \qquad L_2 =$$

L1 - leads to sparse weight matrix, but is not differentiable at 0 L2 - gives lower errors

Andrey Ustyuzhanin

To avoid too complicated models regularisation terms are added to plain loss

 $||W||_{2}^{2}$



Playing with Neural Nets in a browser

<u>Convnet.js</u>:

layer_defs > layer_defs layer_defs layer_defs layer_defs net = new net.make trainer = r momentu





Convolutional Neural Networks

Regular Machine Learning (ML) approach

In Physics:

Describe data by variables => ML => signal vs bckg

In a search engine:

Describe web pages by features => ML => estimation of relevance to a query

This approach works well enough in many areas. Features can be either computed manually or also be estimated by ML (*stacking*).



When it didn't work (good enough)?

- which features help to identify a drum on the picture? or a car? or an animal?
- which features are good to detect a particular phoneme in a soundform?







Image recognition today

can classify objects on images with very high accuracy even if those are images of 500 x 500 and we have only 400 output neurones, that's at least $500 * 500 * 400 = 10^8$ parameters

and that's only linear model!

1. can we reduce the number of parameters?

2. is there something we're missing about image data when treating picture as a vector?





Convolution idea

Convolution is a mathematical operation which describes a rule of how to mix two functions or pieces of information:

- input data and
- the convolution kernel mix together to form

a transformed feature map

Input image



Convolution Kernel

$$\begin{array}{cccc} -1 & -1 & -1 \\ -1 & 8 & -1 \\ -1 & -1 & -1 \end{array}$$

Feature map





Convolution in steps

Calculating convolution by sliding image patches over the entire image.

One image patch (yellow) of the original image (green) is multiplied by the kernel (red numbers in the yellow patch), and its sum is written to one feature map pixel (red cell in convolved feature).







Convolved Feature

Image





Pooling

Motivation:

- Reduce layer size by a factor
- Make NN less sensitive to small image shifts

that a specific feature is in the original input volume (there will be a high the other features.

Example: we have a filter for detecting faces. The exact pixel location of the face is less relevant then the fact that there is a face "somewhere at the top"



- The intuition behind this operation (layer in the network) is that once we know activation value), its exact location is not as important as its relative location to



Regularizations redux

Many parameters \rightarrow overfitting

L1, L2 — regularizations are still used

but this isn't enough!

randomly (dropped)

it prevents co-adaptation of neurons and makes the neural network more stable

Andrey Ustyuzhanin



(a) Standard Neural Net



dropout: at each step of optimization some subset of neuron weights is zeroed



Examples of Neural Network layers

- Input
- Convolution
- Pooling
- Fully connected (FC) layers Basically, a FC layer looks at what high level features most strongly correlate to a particular class and has particular weights so that when you compute the products between the weights and the previous layer, you get the correct probabilities for the different classes.
- Dropout



Yann LeCun Network (LeNet)



LeNet-5 based on CNNs — was developed in 1998, and how it works.

Quiz: 1) what is the convolution size at the 1st step? 2) what is pooling window size at the 2nd step?



Problems with deep architectures



In deep architectures, the gradient varies between parameters drastically $(10^3 - 10^6 \text{ times difference and this is not a limit})$.

This problem usually referred to as 'gradient diminishing'.

Adaptive versions of SGD are used today.



Deep learning original idea (Hinton, 2006)

Deep Belief Network (Hinton et al, 2006) was one of the first ways to train deep structures.

- trained in an unsupervised manner layer-by-layer >
- next layer should "find useful dependencies" in previous layer output
- finally, a classification done using output of last layer (optionally, minor fine-tuning of the network might also required)

Why this is important?

- we can recognize objects and learn patterns without supervision gradient
- backpropagation in the brain??

Practical? No, today deep learning relies mostly on stochastic gradient optimization.







Andrey Ustyuznanin

Elephants

Chairs









Inference of high-level information

Ability of neural networks to extract high-level features can be demonstrated with artistic style transfer.

Image is smoothly modified to preserve the activations on some convolutional layer.











Sketch to picture (texturing)

Online demo to transform a doodle (sketch) into Claude Monet or Vincent Van Gogh:

http://likemo.net/ #sketch_canvas



AlphaGo

too many combinations to check

no reliable way to assess the quality of position

AlphaGo has 2 CNNs:

- 1. network predicts probability to win
- 2. policy network predicts next step of a player

both are given the board (with additional features for each position on the board)






Pitfalls of blind machine learning

CNNs were applied to jet substructure analyses and shown promising results on simulated data, but a study shows model applied to different MC generators images has up to 2x background efficiency.

- thus, learnt dependencies are MC-specific and may turn out to be useless,
- machine learning may use dependencies there are approaches to domain adaptation, that might come handy.













Applications of CNNs

- Convolutional neural networks are topperforming on the image and image-like data:
- image recognition and annotation
- object detection
- image segmentation
- segmentation of neuronal membranes

Next big thing: <u>attention control</u> networks

Andrey Ustyuzhanin



38

Recurrent Neural Networks

Sequences and Recurrent NN

sound, text any information that is time dependent (daily revenues, stocks) prices...)

Inspiration: brain contains information about past and is able to process the information given in sequences.





LSTM: a popular recurrent unit

LSTMs have chain like structure, but the repeating module has a different structure. Instead of having a single neural network layer, there are four, interacting in a very special way.

See <u>post</u> by Christopher Olah for details.



RNN for fun, baby names generation

feed the RNN a bunch of 8000 baby names and ask to generate new ones:

Rudi Levette Berice Lussa Hany Mareanne Chrestina Carissy Marylen Hammine Janye Marlise Jacacrie Hendred Romand Charienna Nenotto Ette Dorane Wallen Marly Darine Salina Elvyn Ersia Maralena Minoria Ellia Charmin Nerille Chelon Walmor Jeryly Stachon Allisa Anatha Cathanie Geetra Alexie Jerin Cassen Herbett Cossie Velen Daurenge **Robester** Shermond Terisa Licia Roselen Ferine Jayn Lusine Charyanne Sales Wallon Martine Merus Jelen Candica Wallin Tel Rachene Tarine Ozila Ketia Shanne Arnande Karella Roselina Alessia Chasty Deland Berther Geamar Jackein Mellisand Sagdy Nenc Lessie Guen Gavi Milea Anneda Margoris Janin Rodelin ZeElyne Janah Ferzina Pey Castina, Killie, Saddie, R



RNN for fun, math text generation

Proof. Omitted.

Lemma 0.1. Let C be a set of the construction.

Let C be a gerber covering. Let F be a quasi-coherent sheaves of O-modules. We have to show that

$$\mathcal{O}_{\mathcal{O}_X} = \mathcal{O}_X(\mathcal{L})$$

Proof. This is an algebraic space with the composition of sheaves \mathcal{F} on $X_{\acute{e}tale}$ we have

 $\mathcal{O}_X(\mathcal{F}) = \{morph_1 \times_{\mathcal{O}_X} (\mathcal{G}, \mathcal{F})\}$

where \mathcal{G} defines an isomorphism $\mathcal{F} \to \mathcal{F}$ of \mathcal{O} -modules.

Lemma 0.2. This is an integer Z is injective.

Proof. See Spaces, Lemma ??.

Lemma 0.3. Let S be a scheme. Let X be a scheme and X is an affine open covering. Let $\mathcal{U} \subset \mathcal{X}$ be a canonical and locally of finite type. Let X be a scheme. Let X be a scheme which is equal to the formal complex.

The following to the construction of the lemma follows.

Let X be a scheme. Let X be a scheme covering. Let

 $b: X \to Y' \to Y \to Y \to Y' \times_X Y \to X.$

be a morphism of algebraic spaces over S and Y.

Proof. Let X be a nonzero scheme of X. Let X be an algebraic space. Let \mathcal{F} be a quasi-coherent sheaf of \mathcal{O}_X -modules. The following are equivalent

(1) \mathcal{F} is an algebraic space over S.

(2) If X is an affine open covering.

Consider a common structure on X and X the functor $\mathcal{O}_X(U)$ which is locally of finite type.





RNN for fun, (Linux) source code generation

```
* Increment the size file of the new incorrect UI_FILTER group information
 * of the size generatively.
 */
static int indicate_policy(void)
 int error;
 if (fd == MARN_EPT) {
    /*
     * The kernel blank will coeld it to userspace.
     */
    if (ss->segment < mem_total)</pre>
      unblock_graph_and_set_blocked();
    else
      ret = 1;
    goto bail;
  segaddr = in_SB(in.addr);
 selector = seg / 16;
  setup_works = true;
  for (i = 0; i < blocks; i++) {</pre>
    seq = buf[i++];
   bpf = bd->bd.next + i * search;
   if (fd) {
      current = blocked;
    ٦
  }
 rw->name = "Getjbbregs";
  bprm self clearl(&iv->version);
  regs->new = blocks[(BPF_STATS << info->historidac)] | PFMR_CLOBATHINC_SECONDS << 12;</pre>
  return segtable;
```



Style in RNNs

Take the brook away show they are He dismissed the idea prison welfare Officer compensat

Position of the pen in the next moment is predicted (<u>online demo</u>).

She looked dosely as she and Huntercombe in being adapted for



Different strategies for prediction



one can predict

a single observation from sequence or map sequence to sequence

Andrey Ustyuzhanin

many to many



many to many



Applications of Recurrent NNs

natural language processing /translation speech recognition

sound record is a sequence!

speech synthesis demand forecasting in retail systems / stock prices forecasting Jet substructure recognition, https:// arxiv.org/abs/1702.00748





Mixing RNN and CNN

Sometimes you need both to analyze dependencies in space and time. For example, in weather forecasting a sequence of photos from satellites can be used to predict weather conditions, i.e. snow/rain (https://yandex.com/company/blog/

///)

18:20	2	Сертолово Песочный Кузьмоловский
18:30	2	Сестрорецк Парголово Бугры Курортный район Мурино Роман
18:40	2	Кронштадт
18:50	2	Ажора Невская Губа Санкт-Петербургос Старая Ломоносов Заневское поселение
19:00		поселение Стрельна свероповское поселен
19:10	2	Кицкое поселение
19:20	2	поселение Русско-Высоцкое Ни Павловск-
19:30	2	Кићенское поселение Малое Верево, Коммунар
19:40	2	левское поселение Гатчина Форносово





Other directions

 \searrow

Autoencoding



Encode

compressed data is treated as a result of an algorithm trained representations can be made stable against different noise



Autoencoding



Take MNIST dataset (left) and compare PCA (middle) and autoencoder with 2 neurons in the middle layer (right).

Interactive playground: <u>https://transcranial.github.io/keras-js/#/</u>





Generative Adversarial Networks



how can network learn to generate an image (of a bedroom)? (https://github.com/Newmu/dcgan_code)



Generative Adversarial Networks

Goal: learn to generate realistic images Solution: make Generator play cooperative game with Discriminator

- Generator (MLP) converts noise (Z) to image (S)
- Discriminator (MLP) distinguishes real images (X) from generated > images (X)

Combined loss function:

$$\mathscr{L}_{adv} = E[log(P(D(I)=0|I\in S))] + E[log(P(D(I)=1|I\in X))]$$

term associated with discriminator perceiving a generated sample as fake

The game-theoretical basis for this framework ensures that if we extend the space of allowed functions that G and D can draw from to the space of all continuous functions, then there exists some G that exactly recovers f: Z -> X



- term associated with discriminator perceiving a real sample as real







Applications. Detecting cosmic rays

- cosmic rays coming from space have huge energies
- we don't know their source
- decay products of those can be detected need to cover large areas with detectors
- Collaboration: CRAYFIS, https://crayfis.io
- Idea: use smartphone camera as a detector
- **Problem:** how to simulate well camera response for the particles?





Transforming adversarial networks









Andrey Ustyuzhanin

We can 'calibrate' the simulation with adversarial networks. That's real data images:

TAN input is not a noise, but output of a simulator. Here is G output by epoch:







More applications to HEP

https://arxiv.org/abs/1705.02355 - M. Paganinin et al, CaloGAN: Simulating 3D High Energy Particle Showers in Multi-Layer Electromagnetic Calorimeters with Generative Adversarial Networks

https://arxiv.org/abs/1701.05927 - L. de Oliveira et al, Learning Particle Physics by Example: Location-Aware Generative Adversarial Networks for Physics Synthesis



Reinforcement Learning

Reinforcement Learning (RL) is the branch of machine learning that is concerned with making sequences of decisions. It assumes that there is an *agent* that is situated in an *environment*.

At each step, the agent takes an *action*, and it receives an observation and reward from the environment.

An RL algorithm seeks to maximize the agent's total reward, given a previously unknown environment, through a learning process that usually involves lots of trial and error.





Platform, Typical Problems

OpenAl - A toolkit for developing and comparing reinforcement learning algorithms. It supports teaching agents everything from walking to playing games like Pong or Go.





Sign in with GitHub

algorithm on Breakout-v0

0

🍘 neale

2 months ago



Best 100-episode average reward was 776.33 ± 78.86. (Breakout-v0 does not have a specified reward threshold at which it's considered solved.)











Why Neural Networks?

Reinforcement learning algorithms are focused on learning a *policy* (function that takes in observations and outputs) actions (e.g., motor torques)). Other algorithms are focused on learning *value* functions, which measure the goodness of states (i.e., the state of the world) and actions.

In other words, the Q-function defines a policy. We can use NN to represent the policies and Q-functions. For example, when playing Atari games, the input to these networks is an image of the screen, and there is a discrete set of actions. As a policy for this task, you can use a convolutional neural network with similar architecture, which instead outputs the probability of each action. (Guo et al., 2014, Schulman et al., 2015). (more details) Andrey Ustyuzhanin





Contro

AgentNet

Deep Reinforcement Learning library for humans

- Playing Atari with Convolutional NN via OpenAl Gym
 - \circ Can switch to any visual game thanks to awesome Gym interface
 - Very simplistic, non-recurrent suffering from atari flickering, etc.
- Deep Recurrent Kung-Fu training with GRUs and actor-critic • Uses the "Playing atari" example with minor changes • Trains via Advantage actor-critic (value+policy-based)
- Simple Deep Recurrent Reinforcement Learning setup Trying to guess the interconnected hidden factors on a synthetic problem setup

Possible HEP application: detector design optimisation.

Needs a decent simulator!





Deep Learning application examples for HEP

Jet flavour identification:

- https://arxiv.org/abs/1407.5675 CNN, Josh Cogan et al;
- https://arxiv.org/abs/1603.09349 DNN for jets, Pierre Baldi et al;
- https://arxiv.org/abs/1701.05927 GAN for jets, Luke de Oliveira et al;
- https://arxiv.org/abs/1702.00748 RNN for jets, Gilles Louppe et al;

Ultimate application:



HEP Feature Engineering down to discovery









Practical Aspects



Popular libraries

Theano - low level symbolic computation, <u>https://github.com/</u> Theano/Theano Tensorflow, https://www.tensorflow.org

- Keras, https://keras.io/
- Tensorboard

Lasagna, https://github.com/Lasagne/Lasagne Torch, https://github.com/torch

PyTorch, http://pytorch.org

If working on big dataset, think of getting a GPU



Potential (popular) caveats

Neural Nets require much more attention to get results

- get enough training data;
- network initialisation (random weights);
- playing with optimisation procedure;
- training sample nature;
- choosing cost function;
- choosing architecture;
- debugging practice/tools.



Helpful resources

Winter school on Deep Learning: <u>https://github.com/</u> yandexdataschool/CSC_deeplearning MLHEP 2016: https://github.com/yandexdataschool/mlhep2016 MLHEP 2017: bit.ly/mlhep2017, Reading UK, 17-23 Jul Advanced Machine Learning Specialisation at Coursera (to be started fall 2017)



Conclusion

NNs are very popular today and are a subject of intensive research differentiable operations with tensors provide a very good basis for defining useful predicting models Structure of a problem can be effectively used in the structure of the model

DL requires much data and many resources that became available recently and also numerous quite simple tricks to train it better NNs are awesome, but not guaranteed to work well on problems without specific structure



Thank you for attention!

anaderi@yandex-team.ru

anaderiRu@twitter



Special Thanks to

Tatiana Likhomanenko Fedor Ratnikov Denis Derkach Maxim Borisyak Mikhail Hushchyn and all YSDA research team for helping crafting these slides



Backup Slides

Andrey Ustyuzhanin

7



References

«Pattern Recognition and Event Reconstruction in Particle Physics Experiments» <u>http://arxiv.org/abs/</u>physics/0402039

«Pattern recognition» http://bit.ly/28KWfct

«Pattern recognition in HEP» http://bit.ly/28LUPSy

«Современные методы обработки данных в физике высоких энергий» http:// <u>www1.jinr.ru/Pepan/</u> <u>v-33-3/v-</u>33-3-11.pdf

«Performance Evaluation of RANSAC Family» http://www.bmva.org/bmvc/2009/Papers/Paper355/Paper355.pdf

https://en.wikipedia.org/wiki/Kalman_filter

Schaffer, Cullen. "A conservation law for generalization performance." Proceedings of the 11th international conference on machine learning. 1994.

Wolpert, David H. "The supervised learning no-free-lunch theorems." Soft computing and industry. Springer London, 2002. 25-

Wolpert, David H., and William G. Macready. "No free lunch theorems for optimisation." IEEE transactions on evolutionary computation 1.1 (1997): 67-82.

Andrey Ustyuzhanin

- Schmidhuber, Jürgen. "Deep learning in neural networks: An overview." Neural networks 61 (2015): 85-117

71

Moar References

Cohen, et al 2015, On the Expressive Power of Deep Learning: A **Tensor Analysis** http://cs.stanford.edu/people/karpathy/convnetjs/ http://cs.nyu.edu/~rostami/presentations/L1_vs_L2.pdf https://stats.stackexchange.com/questions/2691/making-senseof-principal-component-analysis-eigenvectors-eigenvalues http://www.cs.toronto.edu/~graves/handwriting.cgi? text=nec+hercules+contra+plures&style=&bias=0.15&samples=3 A Connection Between Generative Adversarial Networks, Inverse Reinforcement Learning, and Energy-Based Models https:// arxiv.org/pdf/1611.03852.pdf


Problem 1: Data Certification (CMS)

- Traditionally, quality of the data at CERN CMS experiment is determined manually which requires considerable amount of human efforts;
- ML can save some of those efforts:



Results



http://bit.ly/210MLiN

The aim is to minimise the Manual work with low Loss Rate ("good" classified as "bad") and Pollution Rate ("bad" classified as "good");

~80% saving on manual work is feasible for Pollution & Loss rate of 0.5%.

Next steps: adopt technique for 2016 data & run in production







Machine Learning Challenges

Complications

Lumisection representation Feature engineering Continuous quality update

Algorithms:

Supervised learning, binary classification:

Neural Networks, Gradient Boosting

Active Learning



Problem 2: LHCb Particle Identification (PID)

> Problem: identify charged particle associated with a track (multiclass classification problem)

particle types: Electron, Muon, Pion, Kaon, Proton and Ghost;

LHCb detector provides diverse plentiful information, collected by subdetectors: CALO, RICH, Muon and Track observables, his information should be **combined**;

Monte Carlo-simulated samples.



C. Lippmann - 2003





Neural Networks: Stacking and Special



network output

Andrey Ustyuzhanin



representations for subdetectors are concatenated

network output

Models AUCs

	Ghost	Electron	Muon	Pion	Kaon	Proton
baseline	0.9484	0.9854	0.9844	0.9345	0.9147	0.9178
keras DL	0.9632	0.9914	0.9925	0.9587	0.9319	0.9320
XGBoost	0.9609	0.9908	0.9922	0.9568	0.9303	0.9302
special BDT	0.9636	0.9913	0.9926	0.9576	0.9309	0.9310

- > sample
- BDT has similar quality to keras DL
- Yumber of classes
 Yumber of classes

ROC AUC - a generic ML quality metric, deviation is ~10⁻⁴, due to large training/testing



Improving PID with flat models



Ideal world

Information from subdetectors strongly depends on particle momentum (energy), that leads to strong dependency between PID efficiency and momentum. Undesirable for physics analyses.

Andrey Ustyuzhanin



Real world



Flat Model vs Baseline

Uniform boosting suppresses this dependency:

based on gradient boosting approach modified loss-function to have «unflatness» that penalises for «bumps»

https://arxiv.org/pdf/ 1410.4140v1.pdf

100

80

Efficiency





Uniform boosting provides flatness along 4 variables at once









Machine Learning Challenges & Methods

Data representation (particle traces) Model blending/ensembling from different sub-detectors **Metric Selection**

Multiclassification? One-vs-One? One-vs-all? Accuracy? Logloss? ROC-AUC?

Reduce model output dependence on momentum (flatness)

Methods:

Multi-class classification Deep NN Advanced Boosting (altering loss function)

Andrey Ustyuzhanin

http://bit.ly/2l0vvXc



Problem 3: LHCb Topological Trigger



Andrey Ustyuzhanin

[mm]



LHCb Topological Trigger

Generic trigger for decays of beauty and charm hadrons; Part of Software trigger; Inclusive for any B decay with at least 2 charged daughters including missing particles; Look for 2, 3, 4 track combinations in a wide mass range.





Machine Learning Challenges & Methods

Definition of event (variable number of particles) Training subsample selection Training scheme (different decays) Metric selection Real-time demand, quality-speed trade-off

Methods

Binary classification Model blending Feature selection Model speed-up



Online part using Bonsai BDT

Features hashing using bins before training

Converting decision trees to n-dimensional table (lookup table)

Table size is limited in RAM (1Gb), thus count of bins for each features should be small (5 bins for each of 12 features)

Discretisation reduces the quality by ~1%







Trigger optimisation results



http://arxiv.org/abs/1510.00572

Andrey Ustyuzhanin

N-Body trigger Performance Comparison

Problem 4: $\tau \rightarrow \mu \mu^+ \mu^-$

Search for very-very rare decay (10⁻⁴⁰ according to standard model); Current sensitivity of LHCb is about 10⁻¹⁰; Data sample is selected from what has been collected by triggers; etc), so data and results should be treated under certain assumptions.



- Sits on the top of data analysis chain (after tracking, triggers, preselections,



ML-flavoured sub-problem

- Every decay candidate described by set of high-level features; > Classification: differentiate decays containing signal from others
- (background);
- **Simulated** sample of signal, **real** sample for background but, model should not pick simulation-specific information;
- Trained model output should not correlate with mass of mother particle.

Results:

- http://arxiv.org/abs/1409.8548
- https://www.kaggle.com/c/flavours-of-physics
- Data Doping: <u>http://bit.ly/2IJSEzU</u>
- https://github.com/yandexdataschool/hep_ml/





ML Challenges

Constrained classification:

- flatness;
- signal/background vs MC/real-data check. Metric?
- Prefer classifier with higher number of true positive with lowest possible false positive number;
- Chosen metric: constrains + weighted ROC AUC.



Problem, HEP	Experiment	
Particle Identification	LHCb	
MC generation optimization	SHiP	
Tracking	LHCb, SHiP, COMET	
Jet identification	LHCb	
Triggers	CRAYFIS	
Data modelling	CRAYFIS	
Anomaly Detection, data certification	LHCb	
Triggers	LHCb	
Detector optimisation	SHiP	
Andrey Ustyuzhanin	-	

ML methods
DNN, classification, advanced Boosting
GP, model calibration, non-convex optimisation
Tracking, Clustering, real-time
CNN, multi classification
Enhanced Convolutional Neural Nets (CNN)
Generative Adversarial Nets (GAN)
Time Series, Binary classification
Classification, real-time
Surrogate modelling

