Machine learning

Yann Coadou

CPPM Marseille

Ecoles d'été France Excellence 2019 Physics of the two infinities Marseille, 2 July 2019





Outline





- Introduction
- Optimal discrimination
 - Bayes limit
 - Multivariate discriminant
- Machine learning
 - Supervised and unsupervised learning
- Multivariate discriminants
 - Quadratic and linear discriminants
 - Decision trees
 - Neural networks
 - Deep networks

Introduction



Typical problems in HEP

- Classification of objects
 - separate real and fake leptons/jets/etc.
- Signal enhancement relative to background
- Regression: best estimation of a parameter
 - lepton energy, ∉_T value, invariant mass, etc.

Discrimination of signal from background in HEP

- Event level (Higgs searches, ...)
- Cone level (tau-vs-jet reconstruction, ...)
- Lifetime and flavour tagging (b-tagging, ...)
- Track level (particle identification, ...)
- Cell level (energy deposit from hard scatter/pileup/noise, ...)

Introduction



Input information from various sources

- Kinematic variables (masses, momenta, decay angles, ...)
- Event properties (jet multiplicity, sum of charges, brightness . . .)
- Event shape (sphericity, aplanarity, . . .)
- Detector response (silicon hits, dE/dx, Cherenkov angle, shower profiles, muon hits, . . .)

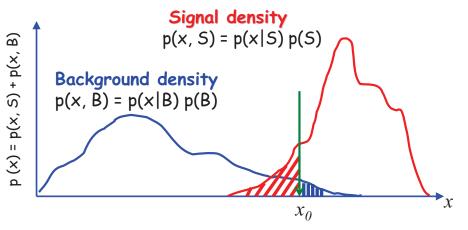
Most data are (highly) multidimensional

- Use dependencies between $x = \{x_1, \dots, x_n\}$ discriminating variables
- Approximate this *n*-dimensional space with a function f(x) capturing the essential features
- f is a multivariate discriminant
- For most of these lectures, use binary classification:
 - ullet an object belongs to one class (e.g. signal) if f(x)>q, where q is some threshold,
 - and to another class (e.g. background) if $f(x) \leq q$

Optimal discrimination: 1-dimension case



• Where to place a cut x_0 on variable x?



• Optimal choice: minimum misclassification cost at decision boundary $x = x_0$

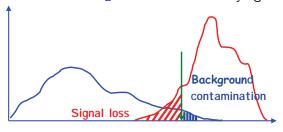
Optimal discrimination: cost of misclassification



$$C(x_0) = C_S \int H(x_0 - x)p(x, S)dx + C_B \int H(x - x_0)p(x, B)dx$$

signal loss background contamination

 C_S = cost of misclassifying signal as background C_B = cost of misclassifying background as signal



- H(x): Heaviside step function
- H(x) = 1 if x > 0, 0 otherwise

Optimal choice: when cost function C is minimum

Optimal discrimination: Bayes discriminant



Minimising the cost

Minimise

$$C(x_0) = C_S \int H(x_0 - x)p(x, S)dx + C_B \int H(x - x_0)p(x, B)dx$$
 with respect to the boundary x_0 :

$$0 = C_S \int \delta(x_0 - x) p(x, S) dx - C_B \int \delta(x - x_0) p(x, B) dx$$

= $C_S p(x_0, S) - C_B p(x_0, B)$

This gives the Bayes discriminant:

$$BD = \frac{C_B}{C_S} = \frac{p(x_0, S)}{p(x_0, B)} = \frac{p(x_0|S)p(S)}{p(x_0|B)p(B)}$$

Probability relationships

- p(A, B) = p(A|B)p(B) = p(B|A)p(A)
- Bayes theorem: p(A|B)p(B) = p(B|A)p(A)
- p(S|x) + p(B|x) = 1

Optimal discrimination: Bayes limit



Generalising to multidimensional problem

• The same holds when x is an *n*-dimensional variable:

$$BD = B \frac{p(S)}{p(B)}$$
 where $B = \frac{p(x|S)}{p(x|B)}$

• B is the Bayes factor, identical to the likelihood ratio when class densities p(x|S) and p(x|B) are independent of unknown parameters

Bayes limit

- p(S|x) = BD/(1+BD) is what should be achieved to minimise cost, achieving classification with the fewest mistakes
- Fixing relative cost of background contamination and signal loss $q = C_B/(C_S + C_B)$, q = p(S|x) defines decision boundary:
 - signal-rich if $p(S|x) \ge q$
 - background-rich if p(S|x) < q
- Any function that approximates conditional class probability p(S|x) with negligible error reaches the Bayes limit

Optimal discrimination: using a discriminant



How to construct p(S|x)?

- k = p(S)/p(B) typically unknown
- Problem: p(S|x) depends on k!
- Solution: it's not a problem...
- Define a multivariate discriminant:

$$D(x) = \frac{s(x)}{s(x) + b(x)} = \frac{p(x|S)}{p(x|S) + p(x|B)}$$

Now:

$$p(S|x) = \frac{D(x)}{D(x) + (1 - D(x))/k}$$

• Cutting on D(x) is equivalent to cutting on p(S|x), implying a corresponding (unknown) cut on p(S|x)

Machine learning: learning from examples



Several types of problems

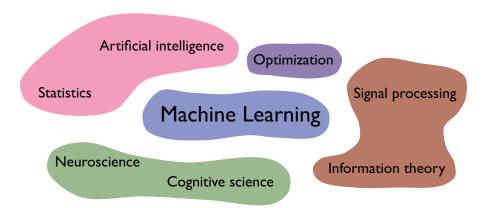
- Classification/decision:
 - signal or background
 - type la supernova or not
 - will pay his/her credit back on time or not
- Regression (mostly ignored in these lectures)
- Clustering (cluster analysis):
 - in exploratory data mining, finding features

Our goal

- ullet Teach a machine to learn the discriminant f(x) using examples from a training dataset
- Be careful to not learn too much the properties of the training sample
 - no need to memorise the training sample
 - instead, interested in getting the right answer for new events
 ⇒ generalisation ability

Machine learning and connected fields





© Balàzs Kégl

Machine learning: (un)supervised learning



Supervised learning

- Training events are labelled: N examples $(x, y)_1, (x, y)_2, \dots, (x, y)_N$ of (discriminating) feature variables x and class labels y
- The learner uses example classes to know how good it is doing

Reinforcement learning

- Instead of labels, some sort of reward system (e.g. game score)
- Goal: maximise future payoff
- May not even "learn" anything from data, but remembers what triggers reward or punishment

Unsupervised learning

- e.g. clustering: find similarities in training sample, without having predefined categories (how Amazon is recommending you books...)
- Discover good internal representation of the input
- Not biased by pre-determined classes ⇒ may discover unexpected features!

Machine learning



Finding the multivariate discriminant y = f(x)

- Given our N examples $(x, y)_1, \dots, (x, y)_N$ we need
 - a function class $\mathbb{F} = \{f(x, w)\}$ (w: parameters to be found)
 - ullet a constraint Q(w) on ${\mathbb F}$
 - a loss or error function L(y, f), encoding what is lost if f is poorly chosen in \mathbb{F} (i.e., f(x, w) far from the desired y = f(x))
- Cannot minimise *L* directly (would depend on the dataset used), but rather its average over a training sample, the empirical risk:

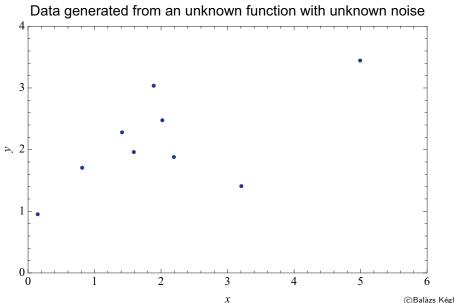
$$R(w) = \frac{1}{N} \sum_{i=1}^{N} L(y_i, f(x_i, w))$$

subject to constraint Q(w), so we minimise the cost function:

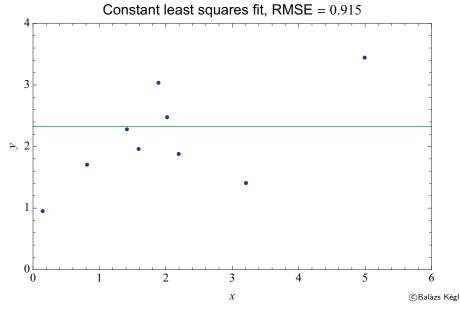
$$C(w) = R(w) + \lambda Q(w)$$

• At the minimum of C(w) we select $f(x, w_*)$, our estimate of y = f(x)

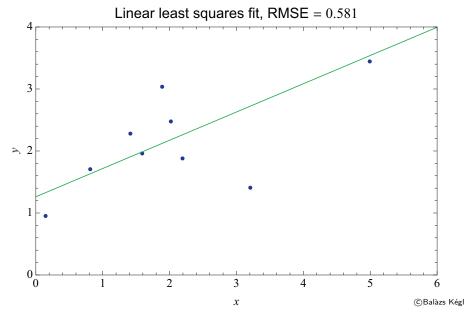




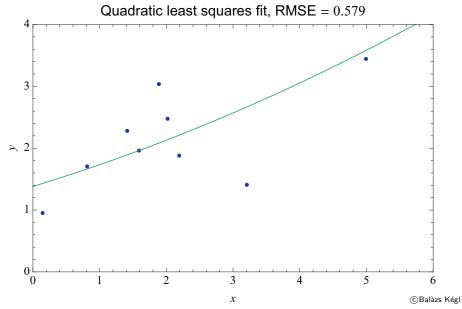




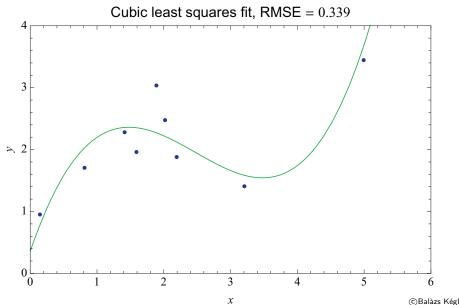




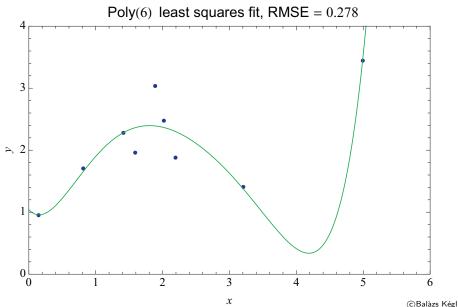




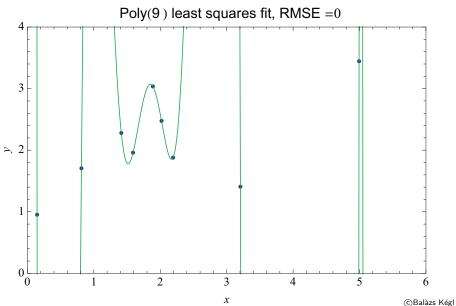














Quality of fit

- Increasing degree of polynomial increases flexibility of function
- Higher degree \Rightarrow can match more features
- If degree = # points, polynomial passes through each point: perfect match!



Quality of fit

- Increasing degree of polynomial increases flexibility of function
- Higher degree \Rightarrow can match more features
- If degree = # points, polynomial passes through each point: perfect match!

Is it meaningful?

- It could be:
 - if there is no noise or uncertainty in the measurement
 - if the true distribution is indeed perfectly described by such a polynomial
- ... not impossible, but not very common...



Quality of fit

- Increasing degree of polynomial increases flexibility of function
- Higher degree \Rightarrow can match more features
- If degree = # points, polynomial passes through each point: perfect match!

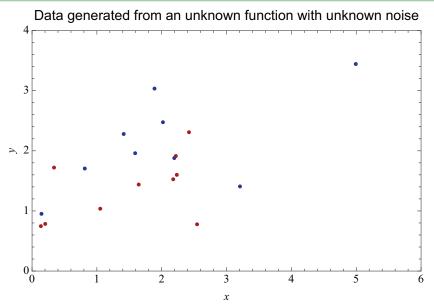
Is it meaningful?

- It could be:
 - if there is no noise or uncertainty in the measurement
 - if the true distribution is indeed perfectly described by such a polynomial
- ... not impossible, but not very common...

Solution: testing sample

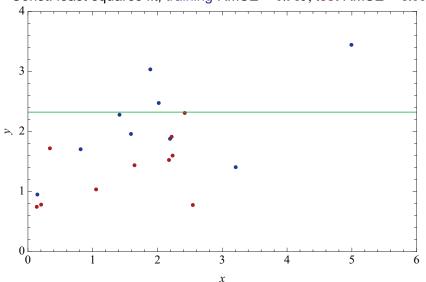
- Use independent sample to validate the result
- Expected: performance will also increase, go through a maximum and decrease again, while it keeps increasing on the training sample





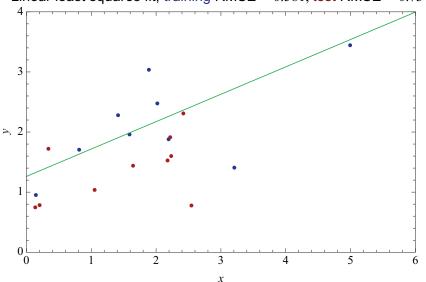






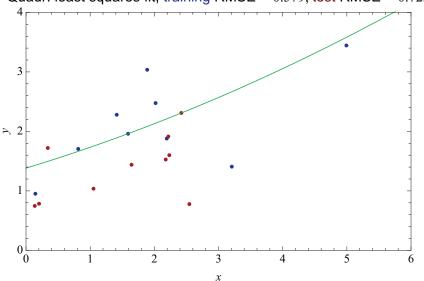








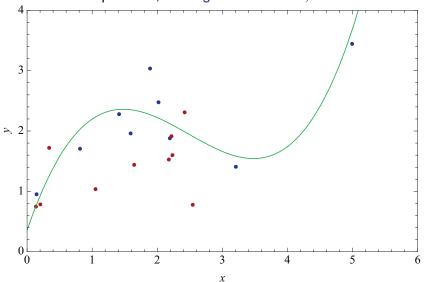
Quadr. least squares fit, training RMSE = 0.579, test RMSE = 0.723



©Balàzs Kégl



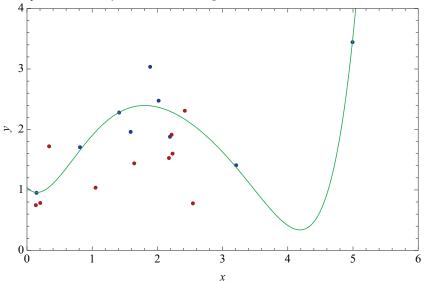




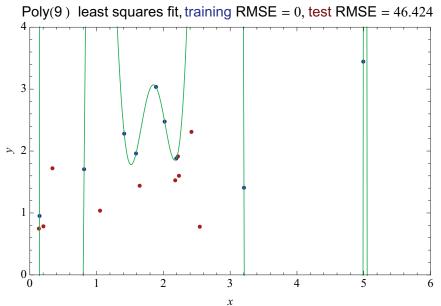
©Balàzs Kégl



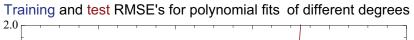


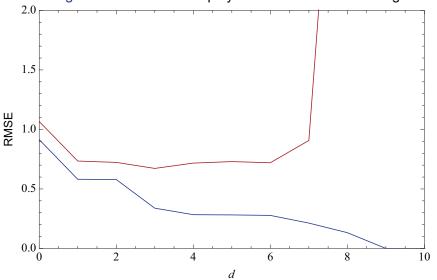














Non-parametric fit

- Minimising the training cost (here, RMSE) does not work if the function class is not fixed in advance (e.g. fix the polynomial degree): complete loss of generalisation capability!
- But if you do not know the correct function class, you should not fix it! Dilemma...

Capacity control and regularisation

- Trade-off between approximation error and estimation error
- Take into account sample size
- Measure (and penalise) complexity
- Use independent test sample
- In practice, no need to correctly guess the function class, but need enough flexibility in your model, balanced with complexity cost

Multivariate discriminants



- Introduction
- Optimal discrimination
 - Bayes limit
 - Multivariate discriminant
- Machine learning
 - Supervised and unsupervised learning
- Multivariate discriminants
 - Quadratic and linear discriminants
 - Decision trees
 - Neural networks
 - Deep networks

Multivariate discriminants



Reminder

• To solve binary classification problem with the fewest number of mistakes, sufficient to compute the multivariate discriminant:

$$D(x) = \frac{s(x)}{s(x) + b(x)}$$

where:

- s(x) = p(x|S) signal density
- b(x) = p(x|B) background density
- Cutting on D(x) is equivalent to cutting on probability p(S|x) that event with x values is of class S

Which approximation to choose?

- Best possible choice: cannot beat Bayes limit (but usually impossible to define)
- No single method can be proven to surpass all others in particular case
- Advisable to try several and use the best one

Quadratic discriminants: Gaussian problem



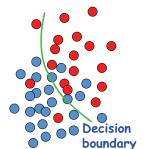
• Suppose densities s(x) and b(x) are multivariate Gaussians:

$$\mathsf{Gaussian}(x|\mu, \Sigma) = \frac{1}{\sqrt{(2\pi)^n |\Sigma|}} \exp^{\left(-\frac{1}{2}(x-\mu)^T \Sigma^{-1}(x-\mu)\right)}$$

with vector of means μ and covariance matrix Σ

• Then Bayes factor B(x) = s(x)/b(x) (or its logarithm) can be expressed explicitly:

$$\ln B(x) = \lambda(x) \equiv \chi^2(\mu_B, \Sigma_B) - \chi^2(\mu_S, \Sigma_S)$$

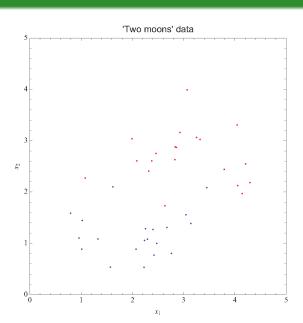


with
$$\chi^2(\mu, \Sigma) = (x - \mu)^T \Sigma^{-1}(x - \mu)$$

- Fixed value of $\lambda(x)$ defines a quadratic hypersurface partitioning the n-dimensional space into signal-rich and background-rich regions
- Optimal separation if s(x) and b(x) are indeed multivariate Gaussians

Quadratic discriminant

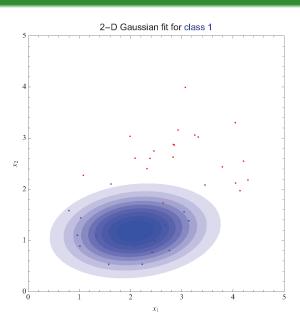




©Balàzs Kégl

Quadratic discriminant

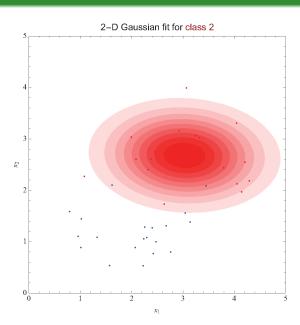




©Balàzs Kégl

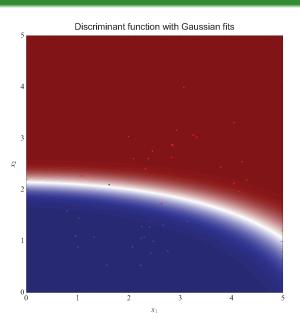
Quadratic discriminant





Quadratic discriminant



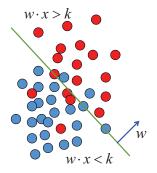


Linear discriminant: Fisher's discriminant



• If in $\lambda(x)$ the same covariance matrix is used for each class (e.g. $\Sigma = \Sigma_S + \Sigma_B$) one gets Fisher's discriminant:

$$\lambda(x) = w \cdot x$$
 with $w \propto \Sigma^{-1}(\mu_S - \mu_B)$



- Optimal linear separation
- Works only if signal and background have different means!
- Optimal classifier (reaches the Bayes limit) for linearly correlated Gaussian-distributed variables

Decision trees



Decision tree origin

 Machine-learning technique, widely used in social sciences. Originally data mining/pattern recognition, then medical diagnostic, insurance/loan screening, etc.



L. Breiman et al., "Classification and Regression Trees" (1984)

Basic principle

- Extend cut-based selection
 - many (most?) events do not have all characteristics of signal or background
 - try not to rule out events failing a particular criterion
- Keep events rejected by one criterion and see whether other criteria could help classify them properly

Binary trees

- Trees can be built with branches splitting into many sub-branches
- In this lecture: mostly binary trees

Tree building algorithm



Start with all events (signal and background) = first (root) node

- sort all events by each variable
- for each variable, find splitting value with best separation between two children
 - mostly signal in one child
 - mostly background in the other
- select variable and splitting value with best separation, produce two branches (nodes)
 - events failing criterion on one side
 - events passing it on the other

Keep splitting

- Now have two new nodes. Repeat algorithm recursively on each node
- Can reuse the same variable
- Iterate until stopping criterion is reached
- Splitting stops: terminal node = leaf



 Consider signal (s_i) and background (b_j) events described by 3 variables: p_T of leading jet, top mass M_t and scalar sum of p_T's of all objects in the event H_T





- Consider signal (s_i) and background (b_j) events described by 3 variables: p_T of leading jet, top mass M_t and scalar sum of p_T 's of all objects in the event H_T
 - sort all events by each variable:
 - $p_T^{s_1} \le p_T^{b_{34}} \le \cdots \le p_T^{b_2} \le p_T^{s_{12}}$
 - $\bullet \ \ H_T^{b_5} \le H_T^{b_3} \le \cdots \le H_T^{s_{67}} \le H_T^{s_{43}}$
 - $M_t^{b_6} \leq M_t^{s_8} \leq \cdots \leq M_t^{s_{12}} \leq M_t^{b_9}$





- Consider signal (s_i) and background (b_j) events described by 3 variables: p_T of leading jet, top mass M_t and scalar sum of p_T 's of all objects in the event H_T
 - sort all events by each variable:

•
$$p_T^{s_1} \le p_T^{b_{34}} \le \cdots \le p_T^{b_2} \le p_T^{s_{12}}$$

•
$$H_T^{b_5} \leq H_T^{b_3} \leq \cdots \leq H_T^{s_{67}} \leq H_T^{s_{43}}$$

•
$$M_t^{b_6} \le M_t^{s_8} \le \cdots \le M_t^{s_{12}} \le M_t^{b_9}$$

- best split (arbitrary unit):
 - $p_T < 56$ GeV, separation = 3
 - $H_T < 242$ GeV, separation = 5
 - $M_t < 105$ GeV, separation = 0.7





- Consider signal (s_i) and background (b_j) events described by 3 variables: p_T of leading jet, top mass M_t and scalar sum of p_T 's of all objects in the event H_T
 - sort all events by each variable:

•
$$p_T^{s_1} \le p_T^{b_{34}} \le \cdots \le p_T^{b_2} \le p_T^{s_{12}}$$

$$\bullet \ \ H_{T}^{b_5} \leq H_{T}^{b_3} \leq \cdots \leq H_{T}^{s_{67}} \leq H_{T}^{s_{43}}$$

•
$$M_t^{b_6} \le M_t^{s_8} \le \cdots \le M_t^{s_{12}} \le M_t^{b_9}$$

- best split (arbitrary unit):
 - $p_T < 56$ GeV, separation = 3
 - $H_T < 242$ GeV, separation = 5
 - $M_t < 105$ GeV, separation = 0.7





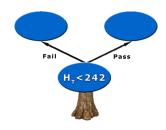
- Consider signal (s_i) and background (b_j) events described by 3 variables: p_T of leading jet, top mass M_t and scalar sum of p_T 's of all objects in the event H_T
 - sort all events by each variable:

•
$$p_T^{s_1} \le p_T^{b_{34}} \le \cdots \le p_T^{b_2} \le p_T^{s_{12}}$$

$$\bullet \ \ H_T^{b_5} \le H_T^{b_3} \le \dots \le H_T^{s_{67}} \le H_T^{s_{43}}$$

•
$$M_t^{b_6} \leq M_t^{s_8} \leq \cdots \leq M_t^{s_{12}} \leq M_t^{b_9}$$

- best split (arbitrary unit):
 - $p_T < 56$ GeV, separation = 3
 - $H_T < 242$ GeV, separation = 5
 - $M_t < 105$ GeV, separation = 0.7
- split events in two branches: pass or fail $H_T < 242 \text{ GeV}$





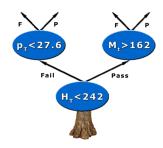
- Consider signal (s_i) and background (b_j) events described by 3 variables: p_T of leading jet, top mass M_t and scalar sum of p_T 's of all objects in the event H_T
 - sort all events by each variable:

•
$$p_T^{s_1} \le p_T^{b_{34}} \le \cdots \le p_T^{b_2} \le p_T^{s_{12}}$$

$$\bullet \ \ H_T^{b_5} \le H_T^{b_3} \le \cdots \le H_T^{s_{67}} \le H_T^{s_{43}}$$

•
$$M_t^{b_6} \leq M_t^{s_8} \leq \cdots \leq M_t^{s_{12}} \leq M_t^{b_9}$$

- best split (arbitrary unit):
 - $p_T < 56$ GeV, separation = 3
 - $H_T < 242$ GeV, separation = 5
 - $M_t < 105$ GeV, separation = 0.7
- split events in two branches: pass or fail $H_T < 242 \text{ GeV}$
- Repeat recursively on each node





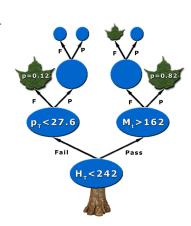
- Consider signal (s_i) and background (b_j) events described by 3 variables: p_T of leading jet, top mass M_t and scalar sum of p_T 's of all objects in the event H_T
 - sort all events by each variable:

•
$$p_T^{s_1} \le p_T^{b_{34}} \le \cdots \le p_T^{b_2} \le p_T^{s_{12}}$$

•
$$H_T^{b_5} \le H_T^{b_3} \le \dots \le H_T^{s_{67}} \le H_T^{s_{43}}$$

•
$$M_t^{b_6} \le M_t^{s_8} \le \cdots \le M_t^{s_{12}} \le M_t^{b_9}$$

- best split (arbitrary unit):
 - $p_T < 56$ GeV, separation = 3
 - H_T < 242 GeV, separation = 5
 - $M_t < 105$ GeV, separation = 0.7
- split events in two branches: pass or fail $H_T < 242 \text{ GeV}$
- Repeat recursively on each node
- Splitting stops: e.g. events with $H_T < 242$ GeV and $M_t > 162$ GeV are signal like (p = 0.82)

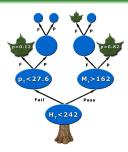


Decision tree output



Run event through tree

- Start from root node
- Apply first best cut
- Go to left or right child node
- Apply best cut for this node
- ...Keep going until...
- Event ends up in leaf



DT Output

- Purity $\left(\frac{s}{s+b}$, with weighted events) of leaf, close to 1 for signal and 0 for background
- or binary answer (discriminant function +1 for signal, -1 or 0 for background) based on purity above/below specified value (e.g. $\frac{1}{2}$) in leaf
- ullet E.g. events with $H_T <$ 242 GeV and $M_t >$ 162 GeV have a DT output of 0.82 or +1

Tree instability: training sample composition

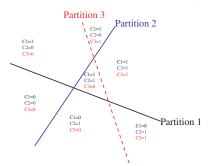


- Small changes in sample can lead to very different tree structures
- Performance on testing events may be as good, or not
- Not optimal to understand data from DT rules
- Does not give confidence in result:
 - DT output distribution discrete by nature
 - granularity related to tree complexity
 - ullet tendency to have spikes at certain purity values (or just two delta functions at ± 1 if not using purity)

Tree (in)stability: distributed representation



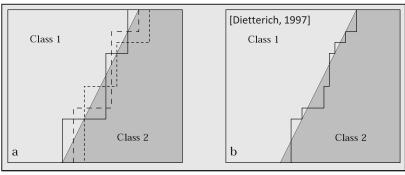
- One tree:
 - one information about event (one leaf)
 - cannot really generalise to variations not covered in training set (at most as many leaves as input size)
- Many trees:
 - distributed representation: number of intersections of leaves exponential in number of trees
 - many leaves contain the event ⇒ richer description of input pattern



Tree (in)stability solution: averaging



Build several trees and average the output



- K-fold cross-validation (good for small samples)
 - divide training sample \mathcal{L} in K subsets of equal size: $\mathcal{L} = \bigcup_{k=1..K} \mathcal{L}_k$
 - Train tree T_k on $\mathcal{L} \mathcal{L}_k$, test on \mathcal{L}_k
 - DT output = $\frac{1}{K} \sum_{k=1..K} T_k$
- Bagging, boosting, random forests, etc.

Boosting: a brief history



First provable algorithm by Schapire (1990)

- Train classifier T_1 on N events
- ullet Train T_2 on new N-sample, half of which misclassified by T_1
- ullet Build T_3 on events where T_1 and T_2 disagree
- Boosted classifier: MajorityVote(T₁, T₂, T₃)

Boosting: a brief history



First provable algorithm by Schapire (1990)

- Train classifier T_1 on N events
- ullet Train T_2 on new N-sample, half of which misclassified by T_1
- Build T_3 on events where T_1 and T_2 disagree
- Boosted classifier: MajorityVote(T₁, T₂, T₃)

Then

- Variation by Freund (1995): boost by majority (combining many learners with fixed error rate)
- Freund&Schapire joined forces: 1st functional model AdaBoost (1996)

Boosting: a brief history



First provable algorithm by Schapire (1990)

- Train classifier T_1 on N events
- ullet Train T_2 on new N-sample, half of which misclassified by T_1
- ullet Build T_3 on events where T_1 and T_2 disagree
- Boosted classifier: MajorityVote(T₁, T₂, T₃)

Then

- Variation by Freund (1995): boost by majority (combining many learners with fixed error rate)
- Freund&Schapire joined forces: 1st functional model AdaBoost (1996)

When it really picked up in HEP

- MiniBooNe compared performance of different boosting algorithms and neural networks for particle ID (2005)
- D0 claimed first evidence for single top quark production (2006)
- CDF copied © (2008). Both used BDT for single top observation

Principles of boosting



What is boosting?

- General method, not limited to decision trees
- Hard to make a very good learner, but easy to make simple, error-prone ones (but still better than random guessing)
- Goal: combine such weak classifiers into a new more stable one, with smaller error

Algorithm

- Training sample \mathbb{T}_k of N events. For i^{th} event:
 - weight w_i^k
 - vector of discriminative variables x_i
 - class label y_i = +1 for signal, -1 for background

Pseudocode:

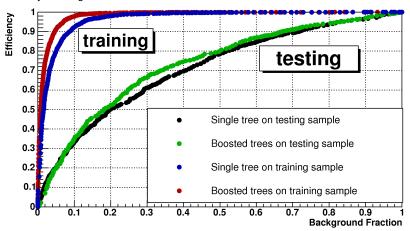
```
Initialise \mathbb{T}_1 for k in 1..N_{tree} train classifier T_k on \mathbb{T}_k assign weight \alpha_k to T_k modify \mathbb{T}_k into \mathbb{T}_{k+1}
```

• Boosted output: $F(T_1, ..., T_{N_{tree}})$

Training and generalisation error



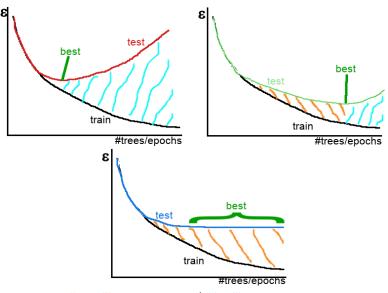
Efficiency vs. background fraction



Clear overtraining, but still better performance after boosting

Overtraining estimation: good or bad?





"good" overtraining / "bad" overtraining

Other averaging techniques



Bagging (Bootstrap aggregating)

- Before building tree T_k take random sample of N events from training sample with replacement
- Train T_k on it
- Events not picked form "out of bag" validation sample
- Applicable to other techniques than DT
 - tends to produce more stable and better classifier

Other averaging techniques



Bagging (Bootstrap aggregating)

- Before building tree T_k take random sample of N events from training sample with replacement
- Train T_k on it
- Events not picked form "out of bag" validation sample
- Applicable to other techniques than DT
 - tends to produce more stable and better classifier

Random forests

- Same as bagging
- In addition, pick random subset of variables to consider for each node split
- Two levels of randomisation, much more stable output
- Often as good as boosting

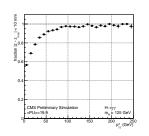
BDT in HEP: CMS H $ightarrow \gamma \gamma$ result

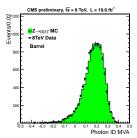


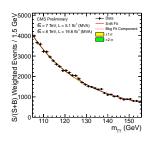
► CMS-PAS-HIG-13-001

Hard to use more BDT in an analysis:

- vertex selected with BDT
- 2nd vertex BDT to estimate probability to be within 1cm of interaction point
- photon ID with BDT
- photon energy corrected with BDT regression
- event-by-event energy uncertainty from another BDT
- several BDT to extract signal in different categories







Neural networks Neuron Human brain • 10^{11} neurons • 10^{14} synapses **Dendrites** • Learning: modifying synapses Electrical Impulses Neurotransmitter Molecules

Receptor

Brief history of artificial neural networks



- 1943: W. McCulloch and W. Pitts explore capabilities of networks of simple neurons
- 1958: F. Rosenblatt introduces perceptron (single neuron with adjustable weights and threshold activation function)
- 1969: M. Minsky and S. Papert prove limitations of perceptron (linear separation only) and (wrongly) conjecture that multi-layered perceptrons have same limitations
 - ⇒ ANN research almost abandoned in 1970s!!!
- 1986: Rumelhart, Hinton and Williams introduce "backward propagation of errors": solves (partially) multi-layered learning
- Next: focus on multilayer perceptron (MLP)

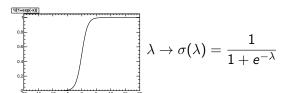
Single neuron

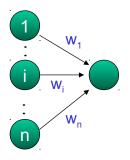


Remember linear separation (Fisher discriminant):

$$\lambda(x) = w \cdot x = \sum_{i=1}^{n} w_i x_i + w_0$$

- Boundary at $\lambda(x) = 0$
- Replace threshold boundary by sigmoid (or tanh):





- $\sigma(\lambda)$ is neuron activity, λ is activation
- ullet Neuron behaviour completely controlled by weights $w=\{w_0,\ldots,w_n\}$
- Training: minimisation of error/loss function (quadratic deviations, entropy [maximum likelihood]), via gradient descent or stochastic approximation

Neural networks



Theorem

Let $\sigma(.)$ be a non-constant, bounded, and monotone-increasing continuous function. Let $\mathcal{C}(I_n)$ denote the space of continuous functions on the n-dimensional hypercube. Then, for any given function $f \in \mathcal{C}(I_n)$ and $\varepsilon > 0$ there exists an integer M and sets of real constants w_j , w_{ij} where $i = 1, \ldots, n$ and $j = 1, \ldots, M$ such that

$$y(x, w) = \sum_{j=1}^{M} w_j \sigma \left(\sum_{i=1}^{n} w_{ij} x_i + w_{0j} \right)$$

is an approximation of f(.), that is $|y(x) - f(x)| < \varepsilon$

Neural networks



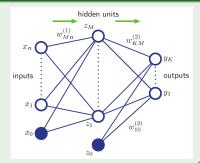
Interpretation

- You can approximate any continuous function to arbitrary precision with a linear combination of sigmoids
- Corollary 1: can approximate any continuous function with neurons!
- Corollary 2: a single hidden layer is enough
- Corollary 3: a linear output neuron is enough

Multilayer perceptron: feedforward network

- Neurons organised in layers
- Output of one layer becomes input to next layer

$$y_k(x, w) = \sum_{j=0}^{M} w_{kj}^{(2)} \underbrace{\sigma\left(\sum_{i=0}^{n} w_{ji}^{(1)} x_i\right)}_{z_j}$$

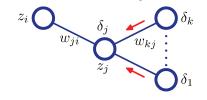


Backpropagation



- Training means minimising error function E(w)
- For single neuron: $\frac{dE}{dw_k} = (y t)x_k$
- One can show that for a network:

$$\frac{dE}{dw_{ji}} = \delta_j z_i$$
, where



$$\delta_k = (y_k - t_k)$$
 for output neurons $\delta_j \propto \sum_k w_{kj} \delta_k$ otherwise

• Hence errors are propagated backwards

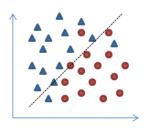
Neural network training

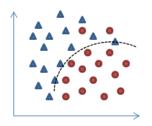


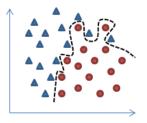
- Minimise error function E(w)
- Gradient descent: $w^{(k+1)} = w^{(k)} \eta \frac{dE^{(k)}}{dw}$
- $\frac{\partial E}{\partial w_j} = \sum_{n=1}^N -(t^{(n)} y^{(n)}) x_j^{(n)}$ with target $t^{(n)}$ (0 or 1), so $t^{(n)} y^{(n)}$ is the error on event n
- All events at once (batch learning):
 - weights updated all at once after processing the entire training sample
 - finds the actual steepest descent
 - takes more time
- or one-by-one (online learning):
 - speeds up learning
 - may avoid local minima with stochastic component in minimisation
 - careful: depends on the order of training events
- One epoch: going through the training data once

Neural network overtraining







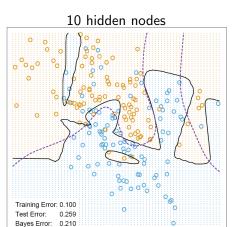


- Diverging weights can cause overfitting
- Mitigate by:
 - early stopping (after a fixed number of epochs)
 - monitoring error on test sample
 - regularisation, introducing a "weight decay" term to penalise large weights, preventing overfitting:

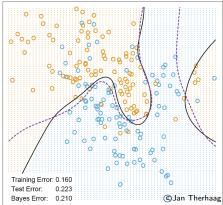
$$\tilde{E}(w) = E(w) + \frac{\alpha}{2} \sum_{i} w_i^2$$

Regularisation









• Much less overfitting, better generalisation properties

Neural networks: Tricks of the trade Pefficient BackProp

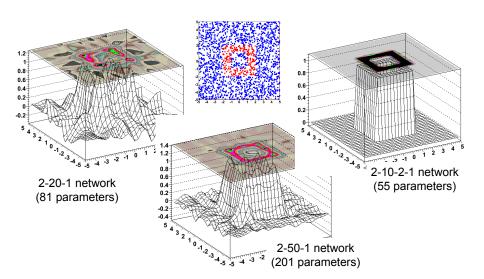




- Preprocess data:
 - if relevant, provide e.g. x/y instead of x and y
 - subtract the mean because the sigmoid derivative becomes negligible very fast (so, input mean close to 0)
 - normalise variances (close to 1)
 - shuffle training sample (order matters in online training)
- Initial random weights should be small to avoid saturation
- Batch/online training: depends on the problem
- Regularise weights to minimise overtraining. May also help select good variables via Automatic Relevance Determination (ARD)
- Make sure the training sample covers the full parameter space
- No rule (not even guestimates) about the number of hidden nodes (unless using constructive algorithm, adding resources as needed)
- A single hidden layer is enough for all purposes, but multiple hidden layers may allow for a solution with fewer parameters

Adding a hidden layer





Deep learning



What is learning?

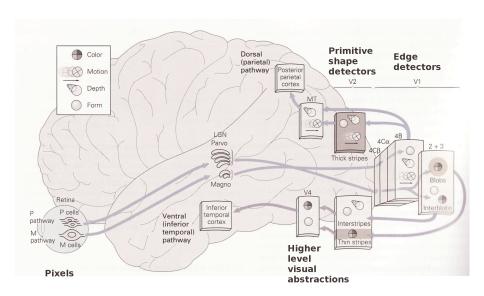
- Ability to learn underlying and previously unknown structure from examples
 - ⇒ capture variations
- ullet Deep learning: have several hidden layers (> 2) in a neural network

Motivation for deep learning

- Just like in the brain!
- Humans organise ideas hierarchically, through composition of simpler ideas
- Heavily unsupervised training, learning simpler tasks first, then combined into more abstract ones
- Learn first order features from raw inputs, then patterns in first order features, then etc.

Deep architecture in the brain





Deep learning in artificial intelligence



Mimicking the brain

- About 1% of neurons active simultaneously in the brain: distributed representation
 - activation of small subset of features, not mutually exclusive
 - more efficient than local representation
 - distributed representations necessary to achieve non-local generalization, exponentially more efficient than 1-of-N enumeration
 - example: integers in 1..N
 - local representation: vector of N bits with single 1 and N-1 zeros
 - distributed representation: vector of log₂ N bits (binary notation), exponentially more compact
- Meaning: information not localised in particular neuron but distributed across them

Deep architecture

- Insufficient depth can hurt
- Learn basic features first, then higher level ones
- Learn good intermediate representations, shared across tasks

Deep learning revolution



Deep networks were unattractive

- One layer is theoretically enough for everything
- Used to perform worse than shallow networks with 1 or 2 hidden layers
- Apparently difficult/impossible to train (using random initial weights and supervised learning with backpropagation)
- Backpropagation issues:
 - requires labelled data (usually scarce and expensive)
 - does not scale well, getting stuck in local minima
 - "vanishing gradient": gradients getting very small further away from output ⇒ early layers do not learn much, can even penalise overall performance

Deep learning revolution



Deep networks were unattractive

- One layer is theoretically enough for everything
- Used to perform worse than shallow networks with 1 or 2 hidden layers
- Apparently difficult/impossible to train (using random initial weights and supervised learning with backpropagation)
- Backpropagation issues:
 - requires labelled data (usually scarce and expensive)
 - does not scale well, getting stuck in local minima
 - "vanishing gradient": gradients getting very small further away from output

 early layers do not learn much, can even penalise overall performance

Breakthroughs around 2006 (Bengio, Hinton, LeCun)

- Try to model structure of input, p(x) instead of p(y|x)
- Can use unlabelled data (a lot of it), with unsupervised training
- Train each layer independently (pre-train and stack)
- New activation functions (e.g. rectified linear unit ReLU)
- Possible thanks to algorithmic innovations, computing resources, data!

Greedy layer-wise pre-training



Algorithm

- Take input information
- Train feature extractor
- Use output as input to training another feature extractor
- Keep adding layers, train each layer separately
- Finalise with a supervised classifier, taking last feature extractor output as input
- All steps above: pre-training
- Fine-tune the whole thing with supervised training (backpropagation)
 - initial weights are those from pre-training

Feature extractors

- Restricted Boltzmann machine (RBM), auto-encoder, sparse auto-encoder, denoising auto-encoder, etc.
- Note: important to not use linear activation functions in hidden layers. Combination of linear functions still linear, so equivalent to single hidden layer

Why does unsupervised training work?



Optimisation hypothesis

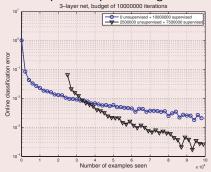
- Training one layer at a time scales well
- Backpropagation from sensible features
- Better local minimum than random initialisation, local search around it

Overfitting/regularisation hypothesis

- More info in inputs than labels
- No need for final discriminant to discover features
- Fine-tuning only at category boundaries

Example

- Stacked denoising auto-encoders
- 10 million handwritten digits
- First 2.5 million used for unsupervised pre-training



 Worse with supervision: eliminates projections of data not useful for local cost but helpful for deep model cost

An example from Google research team • 2011 paper





A "giant" neural network

- At Google they trained a 9-layered NN with 1 billion connections
 - trained on 10 million 200×200 pixel images from YouTube videos
 - on 1000 machines (16000 cores) for 3 days, unsupervised learning
- Sounds big? The human brain has 100 billion (10¹¹) neurons and 100 trillion (10¹⁴) connections...

What it did

- It learned to recognise faces, one of the original goals
- ... but also cat faces (among the most popular things in YouTube videos) and body shapes







Google's research on building high-level features

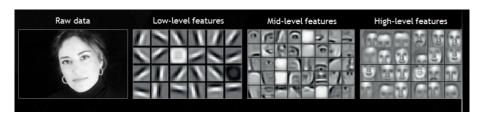




- Features extracted from such images
- Results shown to be robust to
 - colour
 - translation
 - scaling
 - out-of-plane rotation

Learning feature hierarchy





Auto-encoders

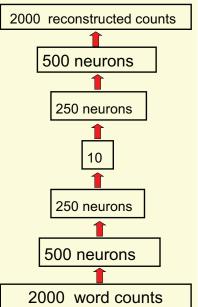


Approximate the identity function

- Build a network whose output is similar to its input
- Sounds trivial? Except if imposing constraints on network (e.g., # of neurons, locally connected network) to discover interesting structures
- Can be viewed as lossy compression of input

Finding similar books

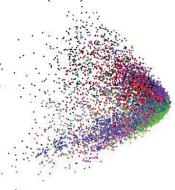
- Get count of 2000 most common words per book
- "Compress" to 10 numbers

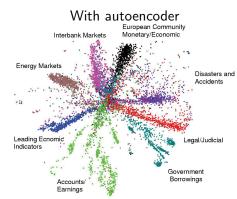


Auto-encoders



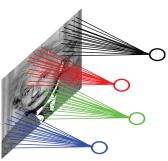
With principle component analysis (PCA)







• Images are stationary: can learn feature in one part and apply it in another





- Images are stationary: can learn feature in one part and apply it in another
- Use e.g. small patch sampled randomly, learn feature, convolve with full image

1 _{×1}	1,0	1,	0	0
O _{×0}	1,	1,0	1	0
O _{×1}	0,0	1,	1	1
0	0	1	1	0
0	1	1	0	0



Image

Convolved Feature



- Images are stationary: can learn feature in one part and apply it in another
- Use e.g. small patch sampled randomly, learn feature, convolve with full image

1	1,	1,0	0,,1	0
0	1 _{×0}	1,	1 _{×0}	0
0	0,,1	1,0	1,	1
0	0	1	1	0
0	1	1	0	0



Image

Convolved Feature



- Images are stationary: can learn feature in one part and apply it in another
- Use e.g. small patch sampled randomly, learn feature, convolve with full image

$\begin{array}{c ccccccccccccccccccccccccccccccccccc$	1	1	1,	0,0	0,,1
0 0 1 1 0	0	1	1,0	1 _{×1}	0,0
	0	0	1,	1,0	1,
0 1 1 0 0	0	0	1	1	0
	0	1	1	0	0



Image

Convolved Feature



- Images are stationary: can learn feature in one part and apply it in another
- Use e.g. small patch sampled randomly, learn feature, convolve with full image

1	1	1	0	0
0,	1,0	1,	1	0
0,	0,1	1,0	1	1
0,	0,×0	1,	1	0
0	1	1	0	0



Image

Convolved Feature



- Images are stationary: can learn feature in one part and apply it in another
- Use e.g. small patch sampled randomly, learn feature, convolve with full image

1	1	1	0	0
0	1 _{×1}	1 _{×0}	1 _{×1}	0
0	0,0	1,	1,0	1
0	0,,1	1 _{×0}	1 _{×1}	0
0	1	1	0	0



Image

Convolved Feature



- Images are stationary: can learn feature in one part and apply it in another
- Use e.g. small patch sampled randomly, learn feature, convolve with full image

0 1 1 1 1	0 _{×1}
0 0 1, 1,	1,0
0 0 1, 1,	0 _{×1}
0 1 1 0	0



Image

Convolved Feature



- Images are stationary: can learn feature in one part and apply it in another
- Use e.g. small patch sampled randomly, learn feature, convolve with full image

1	1	1	0	0
0	1	1	1	0
0,1	0,0	1,	1	1
0,0	0,,1	1,0	1	0
0,1	1,0	1,	0	0

4	3	4
2	4	3
2		

Image

Convolved Feature



- Images are stationary: can learn feature in one part and apply it in another
- Use e.g. small patch sampled randomly, learn feature, convolve with full image

1	1	1	0	0
0	1	1	1	0
0	0,,1	1,0	1,	1
0	0,0	1,	1 _{×0}	0
0	1,	1,0	0,,1	0

4	3
3	
	4 3

Image

Convolved Feature



- Images are stationary: can learn feature in one part and apply it in another
- Use e.g. small patch sampled randomly, learn feature, convolve with full image

1	1	1	0	0
0	1	1	1	0
0	0	1,	1,0	1,
0	0	1,0	1,	0,0
0	1	1,	0,0	0,,1

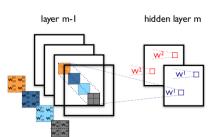
4	3	4
2	4	3
2	3	4

Image

Convolved Feature



- Images are stationary: can learn feature in one part and apply it in another
- Use e.g. small patch sampled randomly, learn feature, convolve with full image
- Build several "feature maps"





WI L

WI D

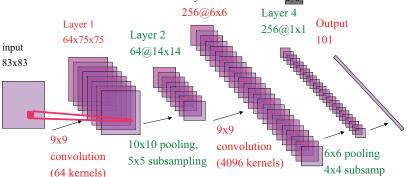
hidden layer m

 Images are stationary: can learn feature in one part and apply it in another

 Use e.g. small patch sampled randomly, learn feature, convolve with full image

Build several "feature maps"

• Stack them with pooling layers_{Layer 3}



layer m-I

Deep learning: looking forward



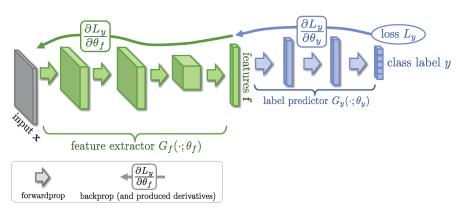
- Very active field of research in machine learning and artificial intelligence
 - not just at universities (Google, Facebook, Microsoft, NVIDIA, etc...)
- Training with curriculum:
 - what humans do over 20 years, or even a lifetime
 - learn different concepts at different times
 - solve easier or smoothed version first, and gradually consider less smoothing
 - exploit previously learned concepts to ease learning of new abstractions
- Influence learning dynamics can have big impact:
 - order and selection of examples matters
 - choose which examples to present first, to guide training and possibly increase learning speed (called shaping in animal training)
- Combination of deep learning and reinforcement learning
 - still in its infancy, but already impressive results
- Domain adaptation and adversarial training
 - e.g. train in parallel network that produces difficult examples
 - learn discrimination (s vs. b) and difference between training and application samples (e.g. Monte Carlo simulation and real data)

Domain adaptation and adversarial training



- Typical training
 - signal and background from simulation

- ▶ http://arxiv.org/abs/1505.07818
- results compared to real data to make measurement
- Requires good data-simulation agreement

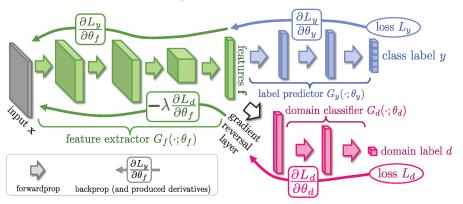


Domain adaptation and adversarial training



- Typical training
 - signal and background from simulation

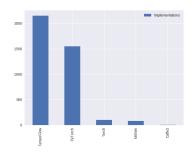
- ► http://arxiv.org/abs/1505.07818
- results compared to real data to make measurement
- Requires good data-simulation agreement
- Possibility to use adversarial training and domain adaptation to account for discrepancies/systematic uncertainties



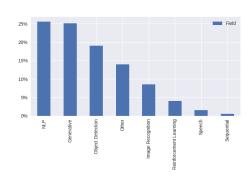
Deep learning: second half of 2018 trends



Most used software



Most activity



https://www.kdnuggets.com/2018/12/deep-learning-major-advances-review.html



ImageNet Large Scale Visual Recognition Challenge

- ImageNet: database with 14 million images and 20k categories
- Used 1000 categories and about 1.3 million manually annotated images

PASCAL



bird



cat



II SVRC



cock













Siamese cat







dalmatian

Egyptian cat



Persian cat

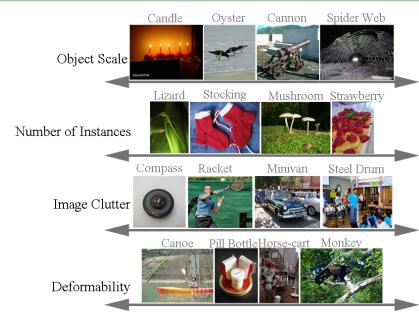






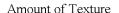
ILSVRC 2014 images





ILSVRC 2014 images







Color Distinctiveness



Shape Distinctiveness



Real-world Size

Low



High

ILSVRC 2014 tasks



Image classification



Ground truth

Steel drum Folding chair Loudspeaker

Accuracy: 1

Scale T-shirt

Steel drum Drumstick Mud turtle

Accuracy: 1

Scale T-shirt Giant panda Drumstick Mud turtle

Accuracy: 0

Single-object localization Steel drum



Ground truth



Accuracy: 1

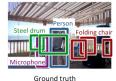


Accuracy: 0



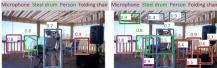
Accuracy: 0

Object detection





AP: 1.0 1.0 1.0 1.0



AP: 0.0 0.5 1.0 0.3



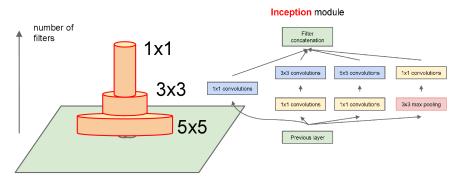
AP: 1.0 0.7 0.5 0.9

ILSVRC 2014 And the winner is...



- Google of course! (first time)
- GoogLeNet:

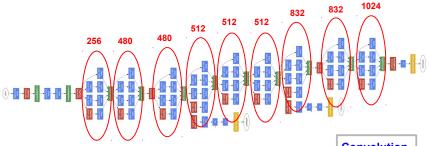
Schematic view



ILSVRC 2014 And the winner is...



- Google of course! (first time)
- GoogLeNet:



9 Inception modules

Network in a network in a network...

Convolution Pooling Softmax Other

ILSVRC 2014 Even GoogLeNet is not perfect!



Classification failure cases



<u>Groundtruth</u>: Police car <u>GoogLeNet</u>:

- laptop
- hair drier
- binocular
- ATM machine
- seat belt

ILSVRC 2010-2016









2010-14: 4.2x reduction

1.7x reduction

1.9x increase

ILSVRC 2015 (same dataset as 2014)

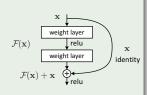
Winner: MSRA (Microsoft Research in Beijing)

• Deep residual networks with > 150 layers

• Classification error: $6.7\% \rightarrow 3.6\%$ (1.9x)

• Localisation error: $26.7\% \rightarrow 9.0\%$ (2.8x)

• Object detection: $43.9\% \rightarrow 62.1\%$ (1.4x)



→ arXiv:1512.03385

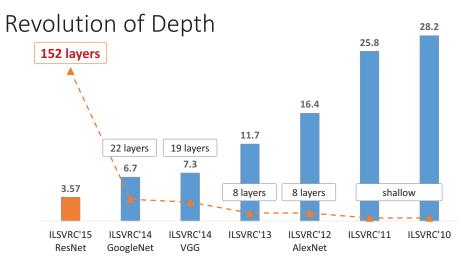
ILSVRC 2016

► http://image-net.org/challenges/LSVRC/2016

Mostly ResNets. Classification: 0.030; localisation: 0.08; detection: 0.66

MSRA @ ILSVRC2015



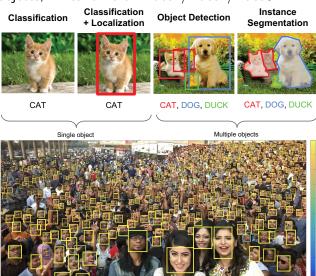


ImageNet Classification top-5 error (%)

Going further



- More and more refinement (segmentation)
- More objects, in real time on video1/video2/video3





- Learning to play 49 different Atari 2600 games
- No knowledge of the goals/rules, just 84x84 pixel frames
- 60 frames per second, 50 million frames (38 days of game experience)
- Deep convolutional network with reinforcement: DQN (deep Q-network)
 - action-value function $Q^*(s,a) = \max_{\pi} \mathbb{E}[r_t + \gamma r_{t+1} + \gamma^2 r_{t+2} + \dots | s_t = s, a_t = a, \pi]$
 - maximum sum of rewards r_t discounted by γ at each timestep t, achievable by a behaviour policy $\pi = P(a|s)$, after making observation s and taking action a
- Tricks for scalability and performance:
 - experience replay (use past frames)
 - separate network to generate learning targets (iterative update of Q)
- Outperforms all previous algorithms, and professional human player on most games

Google DeepMind: training&performance



Algorithm 1: deep O-learning with experience replay.

Initialize replay memory D to capacity N

Initialize action-value function Q with random weights θ Initialize target action-value function \hat{Q} with weights $\theta^- = \theta$

For episode = 1, M do

Initialize sequence $s_1 = \{x_1\}$ and preprocessed sequence $\phi_1 = \phi(s_1)$ For t = 1.T do

With probability ε select a random action a_t

otherwise select $a_t = \operatorname{argmax}_a Q(\phi(s_t), a; \theta)$

Execute action a_r in emulator and observe reward r_r and image x_{r+1} Set $s_{t+1} = s_t, a_t, x_{t+1}$ and preprocess $\phi_{t+1} = \phi(s_{t+1})$

Store transition $(\phi_t, a_t, r_t, \phi_{t+1})$ in D

Sample random minibatch of transitions $(\phi_i, a_j, r_j, \phi_{j+1})$ from D

$$Set y_{j} = \begin{cases} r_{j} & \text{if episode terminates at step } j+1 \\ r_{j} + \gamma \max_{a'} \hat{Q}(\phi_{j+1}, a'; \theta^{-}) & \text{otherwise} \end{cases}$$

Perform a gradient descent step on $(y_j - Q(\phi_j, a_j; \theta))^2$ with respect to the network parameters θ

Every C steps reset $\hat{O} = O$

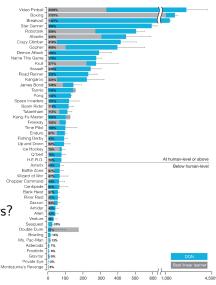
End For End For

• What about Breakout or Space invaders?







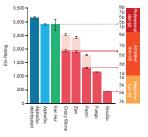




- Game of Go considered very challenging for AI
- Board games: can be solved with search tree of b^d possible sequences of moves (b = breadth [number of legal moves], d = depth [length of game])
- Chess: $b \approx 35$, $d \approx 80 \rightarrow \text{go}$: $b \approx 250$, $d \approx 150$
- Reduction:
 - of depth by position evaluation (replace subtree by approximation that predicts outcome)
 - of breadth by sampling actions from probability distribution (policy p(a|s)) over possible moves a in position s
- ullet 19 imes 19 image, represented by CNN
- Supervised learning policy network from expert human moves, reinforcement learning policy network on self-play (adjusts policy towards winning the game), value network that predicts winner of games in self-play.



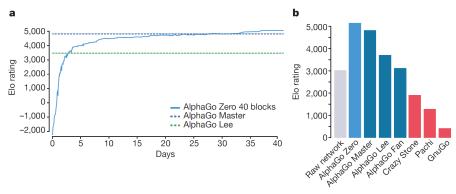
- AlphaGo: 40 search threads, simulations on 48 CPUs, policy and value networks on 8 GPUs. Distributed AlphaGo: 1020 CPUs, 176 GPUs
- AlphaGo won 494/495 games against other programs (and still 77% against Crazy Stone with four handicap stones)
- Fan Hui: 2013/14/15 European champion
- Distributed AlphaGo won 5–0
- AlphaGo evaluated thousands of times fewer positions than Deep Blue (first chess computer to bit human world champion) ⇒ better position selection (policy network) and better evaluation (value network)



- Then played Lee Sedol (top Go play in the world over last decade) in March 2016 ⇒ won 4–1. AlphaGo given honorary professional ninth dan, considered to have "reach a level 'close to the territory of divinity'"
- Ke Jie (Chinese world #1): "Bring it on!". May 2017: 3-0 win for AlphaGo. New comment: "I feel like his game is more and more like the 'Go god'. Really, it is brilliant"

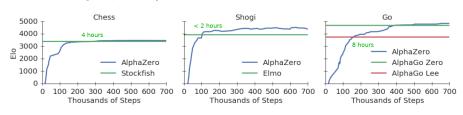


- Learn from scratch, just from the rules and random moves
- Reinforcement learning from self-play, no human data/guidance
- Combined policy and value networks
- 4.9 million self-play games
- Beats AlphaGo Lee (several months of training) after just 36 hours
- Single machine with four TPU





- Same philosophy as AlphaGo Zero, applied to chess, shogi and go
- Changes:
 - not just win/loss, but also draw or other outcomes
 - no additional training data from game symmetries
 - using always the latest network to generate self-play games rather than best one
 - tree search: 80k/70M for chess AlphaZero/Stockfish, 40k/35M for shogi AlphaZero/Elmo

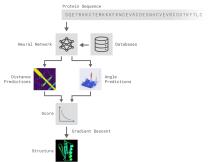


DeepMind AlphaFold





- Trying to tackle scientific problem
- Goal: predict 3D structure of protein based solely on genetic sequence
- Using DNN to predict
 - distances between pairs of amino acids
 - angles between chemical bonds
- Then search DB to find matching existing substructures
- Also train a generative NN to invent new fragments
- Achieved best prediction ever



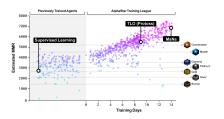


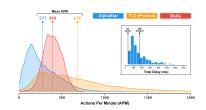
DeepMind AlphaStar





- Mastering real-time strategy game StarCraft II
- Challenges in game theory (no single best strategy), imperfect information (hidden parts of game), long term planning, real time (continuous flow of actions), large action space (many units/buidings)
- Using DNN trained
 - directly on raw data games
 - supervised learning on human games
 - reinforcement learning (continuous league)
- DNN output: list of actions
- Trained for 14 days; each agent: up to 200 years of real-time play
- Runs on single desktop GPU
- Defeated 5–0 one of best pro-players





Deep networks: new results all the time



- Playing poker
 - Libratus (Al developed by Carnegie Mellon University) defeated four of the world's best professional poker players (Jan 2017)
 - After 120,000 hands of Heads-up, No-Limit Texas Hold'em, led the pros by a collective \$1,766,250 in chips
 - Learnt to bluff, and win with incomplete information and opponents' misinformation
- Lip reading ► arXiv:1611.05358 [cs.CV]
 - human professional: deciphers less than 25% of spoken words
 - CNN+LSTM trained on television news programs: 50%

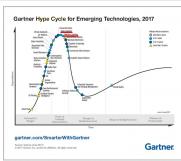


- left: correctly classified image
- middle: difference between left image and adversarial image (x10)
- right: adversarial image, classified as ostrich

Hype cycle

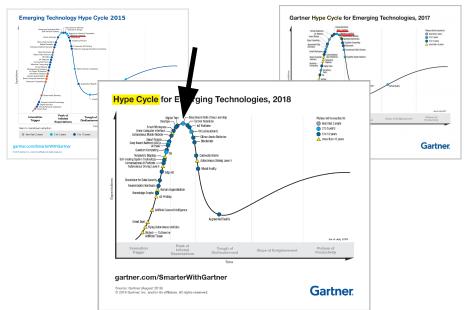




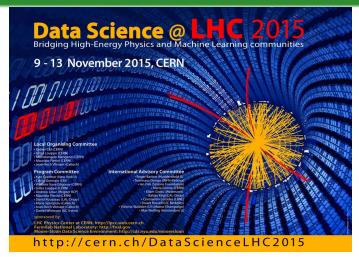


Hype cycle



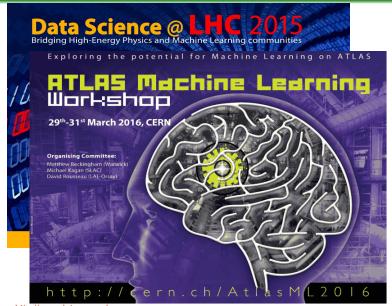






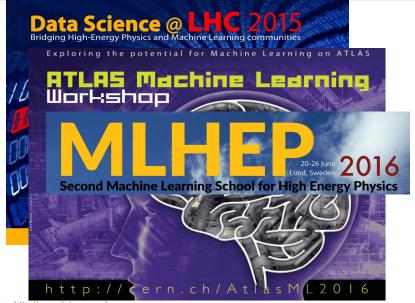
http://opendata.cern.ch





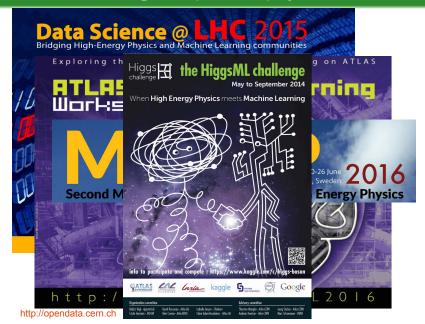
http://opendata.cern.ch



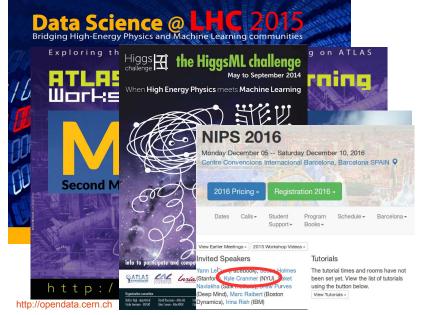


http://opendata.cern.ch









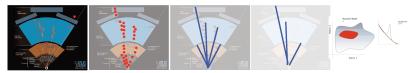




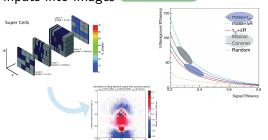


Going to lower level features → arXiv:1410.3469

Raw	Sparsified	Reco	Select	Physics	Ana
1e7	1e4	100-ish*	50	10	1

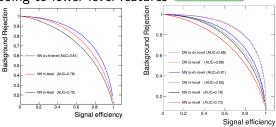


• Transforming inputs into images

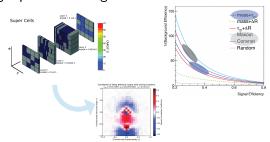




Going to lower level features → arXiv:141

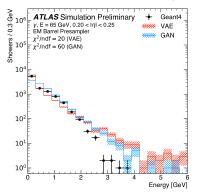


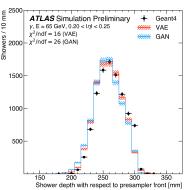
Transforming inputs into images → arXiv:1511.05





- Generative adversarial networks (GAN) ► ATLAS PUB note ATL-SOFT-PUB-2018-001
- Attempts to decrease CPU cost of simulation
 - limiting factor in many analyses already
 - not enough simulated events
- Replace "full simulation" by objects generated automatically by GAN or variational auto-encoders (VAE)





Conclusion



 When trying to achieve optimal discrimination one can try to approximate

$$D(x) = \frac{s(x)}{s(x) + b(x)}$$

- Many techniques and tools exist to achieve this
- (Un)fortunately, no one method can be shown to outperform the others in all cases.
- One should try several and pick the best one for any given problem
- Latest machine learning algorithms (e.g. deep networks) require enormous hyperparameter space optimisation...
- Machine learning and multivariate techniques are at work in your everyday life without your knowning and can easily outsmart you for many tasks

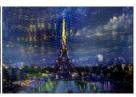
Deep networks and art



• Learning a style ▶ arXiv:1508.06576 [cs.CV]







Computer dreams Google original

Face Style ► http://facestyle.org





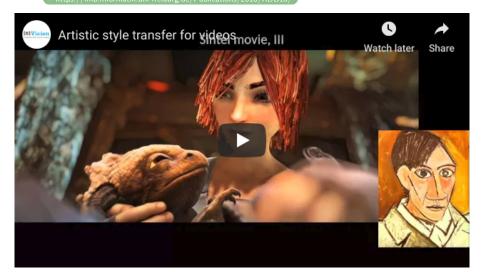


Deep networks and art



Artistic style transfer for videos

https://lmb.informatik.uni_freiburg_de/Publications/2018/RDB18/



References I





T. Hastie, R. Tibshirani and J. Friedman, The Elements of Statistical Learning: Data Mining, Inference and Prediction, Springer-Verlag, New York, 2nd Edition, 2009



C.M. Bishop, Pattern Recognition and Machine Learning, Springer, New York, 2007



M. Minsky and S. Papert, "Perceptrons", M.I.T. Press, Cambridge, Mass., 1969



W.S. McCulloch & W. Pitts, "A logical calculus of the ideas immanent in nervous activity", Bulletin of Mathematical Biophysics, 5, 115-137, 1943



F. Rosenblatt, "The Perceptron: A Probabilistic Model for Information Storage & Organization in the Brain", Psychological Review, 65, pp. 386-408, 1958



D.E.Rumelhart et al., "Learning representations by back-propagating errors", *Nature* vol. 323, p. 533, 1986



K. Hornik et al., "Multilayer Feedforward Networks are Universal Approximators", Neural Networks, Vol. 2, pp 359-366, 1989



Y. LeCun, L. Bottou, G. Orr and K. Muller, "Efficient BackProp", in Neural Networks: Tricks of the trade, Orr, G. and Muller K. (Eds), Springer, 1998

References II





Q.V. Le et al., "Building High-level Features Using Large Scale Unsupervised Learning", in Proceedings of the 29th International Conference on Machine Learning, Edinburgh, Scotland, UK, 2012 http://research.google.com/pubs/pub38115.html



G.E. Hinton, S. Osindero and Y. Teh, "A fast learning algorithm for deep belief nets", Neural Computation 18:1527-1554, 2006



Y. Bengio, P. Lamblin, D. Popovici and H. Larochelle, "Greedy Layer-Wise Training of Deep Networks", in Advances in Neural Information Processing Systems 19 (NIPS'06), pages 153-160, MIT Press 2007



M.A. Ranzato, C. Poultney, S. Chopra and Y. LeCun, in J. Platt et al., "Efficient Learning of Sparse Representations with an Energy-Based Model", in Advances in Neural Information Processing Systems 19 (NIPS'06), pages 1137–1144, MIT Press, 2007



Y. Bengio, "Learning deep architectures for Al", Foundations and Trends in Machine Learning, Vol. 2, No. 1 (2009) 1–127. Also book at Now Publishers http://www.iro.umontreal.ca/lisa/publications2/index.php/publications/show/239



I. Goodfellow, Y. Bengio and A. Courville, "Deep Learning", MIT Press (2016)