

DR ADRIAN BEVAN

PRACTICAL MACHINE LEARNING CLASSIFICATION PROBLEMS WITH NEURAL NETWORKS

LECTURE PLAN

- Classification
- Examples
 - MNIST
 - Introduce the softmax activation function for multiclass output problems.
 - > Particle physics example: Spin-parity assignment of the Higgs boson using $H\to \tau^+\tau^-$

Summary

QMUL Summer School:https://www.qmul.ac.uk/summer-school/Practical Machine Learning QMplus Page:https://qmplus.qmul.ac.uk/course/view.php?id=10006

A. Bevan

CLASSIFICATION

From the introductory NNs Lecture





CLASSIFICATION

- The final perceptron in a network can be used to assign a type to a data example.
 - e.g. consider a binary activation function that gives an all or nothing response with a data sample that contains two types of event (signal and background)
 - Let nothing correspond to one type: background
 - Let all correspond to the other: signal
- Multi-class output is discussed shortly.



y_{0.5}

EXAMPLES: MNIST

- For more than two classification output types we need to have N_{type} perceptrons in the final output layer.
 - Each output perceptron has an all or nothing response that classifies if a training example is classified as that type or not.
 - e.g. the numbers 1, 2, 3, ... 9, 0 [MNIST example] $\int_{0}^{1} \int_{0}^{1} \int_{0}^{1}$
- If we have a complete set of possible outcomes then we can use this constraint to reduce the number of perceptrons to N_{type}.
 - Assumes that the default classification for one category is given by an example not being classified as any of the others.



EXAMPLES: MNIST

- 60000 training examples
- 10000 test examples
- These are greyscale images (one number required to represent each pixel)
 - Renormalise [0, 255] on to [-1, 1] or [0, 1] for processing*.
- Each image corresponds to a 28x28 pixel array of data.
 - For an MLP this translates to 784 features.

* Depends on which activation function is being used. <u>http://yann.lecun.com/exdb/mnist/</u>



EXAMPLES: MNIST

- ▶ The N_{type} = 10 perceptrons are used to make the following decisions:
 - The number 1 vs not the number 1
 - The number 2 vs not the number 2
 - The number 3 vs not the number 3
 - The number 4 vs not the number 4
 - The number 5 vs not the number 5
 - The number 6 vs not the number 6
 - The number 7 vs not the number 7
 - The number 8 vs not the number 8
 - The number 9 vs not the number 9
 - The number 0 vs not the number 0

For those with a statistical background, this is like a null hypothesis and an alternative hypothesis.

The null hypothesis provides a specific response/expectation.

The alternative hypothesis is the complement of the null.

In this context you classify an example as a specific type, or you provide a decision that it is not that type.

We will see more of the MNIST data when talking about convolutional neural networks.



EXAMPLES: MNIST

An alternative representation is to use a softmax activation function to encode the 10 outputs in a single function.

$$f_j(x) = \frac{e^{w_j^T x}}{\sum_{i=1}^N e^{w_i^T x}}$$

i is the index for the output classification type

The score for the ith output is normalised by the sum of outputs.



http://yann.lecun.com/exdb/mnist/

EXAMPLES: MNIST

- LeCun's website lists a number of complicated ways to train a neural network to solve this problem.
- Recent advances in computing have allowed the use of GPU's have meant that MLPs have been applied to the MNIST data, and have produced good results: error rate of 0.35% (Ciresan et al [1]).

ID	architecture	test error for	best test	simulation	weights
	(number of neurons in each layer)	best validation [%]	error [%]	time [h]	[milions]
1	1000, 500, 10	0.49	0.44	23.4	1.34
2	1500, 1000, 500, 10	0.46	0.40	44.2	3.26
3	2000, 1500, 1000, 500, 10	0.41	0.39	66.7	6.69
4	2500, 2000, 1500, 1000, 500, 10	0.35	0.32	114.5	12.11
5	9 × 1000, 10	0.44	0.43	107.7	8.86

[1] <u>Neural Computation, Volume 22, Number 12, December 2010</u> <u>http://yann.lecun.com/exdb/mnist/</u>



EXAMPLES: MNIST

> 35 training examples were mis-classified by the best NN architecture

1 ² 17	1 ¹ 71	q 8 98	ී 9 5 9	9 79	\$ 5 35	°€ 2 3
6 9 4 9	3 5 35	9 ⁴ 9 7	4 9 49	4 ⁴ 9 4	P ² 0 2	5 35
6 ل 16	9 ⁴ 94	b 0 6 0	6 ⁶	४ 6 86	1 ¹ 79	
é 49	0 5 0	5 5 3 5	? 8 98	9 79	77 17	L 1 6 1
27	8-8	$\boldsymbol{\mathcal{F}}^2$	16	65	4 ⁴	Ø
27	58	78	16	65	94	60





- Classification is used in a variety of particle physics scenarios:
 - Make trigger decisions [all or nothing selection]
 - Identify types of particle [provide labels for events to map against charged particle type]
 - Make perform a hypothesis test on some data (e.g. spinparity analysis of the Higgs boson)
- First we introduce a small amount of particle physics to set the context.



EXAMPLES: SPIN-PARITY OF THE HIGGS BOSON

- Atoms are made up of a nucleus that is surrounded by one or more electrons.
- The nucleus contains neutrons and protons - each of these is made of up and down quarks.
- The force carriers play a role in:
 - g: binding nuclear material;
 - γ: binding electrons to atoms.





- Antiparticles exist for these particles*.
- Combinations of: $q\overline{q}$ qqq
- can form new particles that live for a short period of time.
 - Lifetime depends on mass, but is typically 10⁻¹²-10⁻¹⁵ s.





* The photon, γ, is its own antiparticle.

e

EXAMPLES: SPIN-PARITY OF THE HIGGS BOSON

- We build large instruments to detect charged particles.
- There are sets of similar objects to be identified and we can use machine learning and AI to assist in this task.
- e.g. consider charged particles that we reconstruct as "tracks" (c.g. vapour trail indicating the passage of a plane).

Schematic (not to scale) of the penetrating power of charged particles passing through material.





Lets consider the ATLAS detector at CERN's Large Hadron Collider:



Each part of the detector system provides information (electronic pulses) every 25ns.



EXAMPLES: SPIN-PARITY OF THE HIGGS BOSON

Protons are brought together in the heart of the ATLAS detector at the "interaction point".



Many things can happen in the pp collision; we are only interested in very rare events and use pre-selection to identify them for further analysis.



EXAMPLES: SPIN-PARITY OF THE HIGGS BOSON



> This is a $H \rightarrow \tau^+ \tau^-$ candidate event (example)



Context

17

- The Standard Model of particle physics predicts one Higgs boson, and that is a scalar particle.
 - Scalar means that the spin quantum number is zero and that the wave function of the particle ψ (related to the probability distribution) is an even function.
 - \blacktriangleright ψ is a complex number and the probability is given by $P=\psi\psi^*=|\psi|^2$



Contex

18

- > Spin and parity are quantum numbers associated with fundamental particles.
 - Spin is a degree of freedom, which takes the value J=0 for the Standard Model Higgs boson.
 - > Parity is a degree of freedom related to the wave function of the particle;
 - P = +1: Even Parity, means the wave function is even.
 - ▶ P = -1: Odd Parity, means the wave function is odd.
- One question that needs to be addressed is:
 - ▶ What is J^P for the "Higgs boson" discovered at CERN in 2012?
 - Is it a scalar ($J^P = 0^+$), pseudo scalar ($J^P = 0^-$), or some mixture of states?
- To think in terms of hypothesis testing we can test if the Higgs boson is a scalar particle or not; where the not option includes the pseudo scalar and mixture options.



EXAMPLES: SPIN-PARITY OF THE HIGGS BOSON

- > The decay used for this classification problem is $H \rightarrow \tau^+ \tau^-$
- We will encounter data for this channel later in the course.
 - The τ lepton can decay into pairs of quarks, $q\overline{q}$ leptons (e, μ) and neutrinos.
 - The decay to a final state including a ρ or a₁ meson can be used to distinguish between the 0⁺ and 0⁻ or mixture hypotheses.
 - The feature space used for the problem is reconstructed from 2 and 3 body final states of each T decay used to reconstruct the Higgs candidate.



For background, **T** leptons can decay in many ways

Modes with one charged particle						Modes with three charged particles				
Γ_1	particle $^- \geq 0$ neutrals $\ \geq 0 {\cal K}^0 u_ au$		(85.35 ± 0.07) %	S=1.3	Г ₅₆	h	$^-h^-h^+ \ge 0$ neutrals $\ \ge 0 {\cal K}^0_L u_ au$		(15.20 ± 0.08)%	S=1.3
	("1-prong")				Γ ₅₇		$h^- h^- h^+ \geq 0$ neutrals $ u_ au$		(14.57 ± 0.07) %	S=1.3
Γ2	particle $^- \geq 0$ neutrals $\geq 0 {\cal K}^0_L u_ au$		(84.71 ± 0.08) %	S=1.3			(ex. $K^0_S \rightarrow \pi^+\pi^-$)			
Γ ₃	$\mu^- \overline{\nu}_\mu \nu_\tau$	[a]	(17.41 ± 0.04) %	S=1.1			("3-prong")			
Γ4	$\mu^+ \overline{ u}_\mu u_ au \gamma$	[<i>b</i>]	(3.6 \pm 0.4) $ imes$ 10 $^{-3}$		Г ₅₈		$h^- h^- h^+ u_ au$		(9.80 ± 0.07)%	S=1.2
Γ ₅	$e^-\overline{\nu}_e \nu_{\tau}$	[a]	(17.83 ±0.04)%		Γ ₅₉		$h^- h^- h^+ u_ au$ (ex. K^0)		(9.46 ± 0.06)%	S=1.2
Г ₆	$e^-\overline{\nu}_e \nu_{\tau} \gamma$	[b]	(1.75 ±0.18) %		Г ₆₀		$h^- h^- h^+ u_{ au} (ext{ex.} extsf{K}^0, \omega)$		(9.42 ± 0.06)%	S=1.2
Γ ₇	$h^- \ge 0 K_L^0 u_ au$		(12.06 ±0.06) %	S=1.2	Г ₆₁		$\pi^-\pi^+\pi^-\nu_{\tau}$		(9.31 ± 0.06)%	S=1.2
Γ ₈	$h^- \nu_{\tau}$		$(11.53 \pm 0.06)\%$	S=1.2	Г ₆₂	a 1	$\pi^{-}\pi^{+}\pi^{-}\nu_{\tau}(\text{ex.}K^{0})$		(9.02 ± 0.06) %	S=1.1
Γĝ	$\pi^+ u_{ au}$	[a]	(10.83 ±0.06)%	S=1.2	Г ₆₃		$\pi^{-}\pi^{+}\pi^{-}\nu_{\tau}(\text{ex.}K^{0}),$		< 2.4 %	CL=95%
Γ ₁₀	$K^- \nu_{\tau}$	[a]	$(7.00 \pm 0.10) \times 10^{-3}$	S=1.1	-		_ non-axial vector			
Γ_{11}	$h^- \geq 1$ neutrals $ u_ au$		(37.10 ±0.10)%	S=1.2	I 64		$\pi \pi' \pi \nu_{\tau} (\text{ex.} K^\circ, \omega)$	[a]	$(8.99 \pm 0.06)\%$	S=1.1
Γ ₁₂	$h^- \geq 1\pi^0 \nu_{ au}$ (ex. \mathcal{K}^0)		(36.58 ±0.10)%	S=1.2	۱ ₆₅		h h $h' \ge 1$ neutrals ν_{τ}		$(5.39 \pm 0.07)\%$	S=1.2
$\Gamma_{13}^{}$	$h^- \pi^0 \nu_{\tau}$		(25.95 ±0.09)%	S=1.1	۱ ₆₆		$h h h' \ge 1\pi^{\circ}\nu_{\tau}$ (ex. K°)		$(5.09 \pm 0.06)\%$	S=1.2
Γ_{14}	$\rho \pi^- \pi^0 \nu_{\tau}$	[a]	(25.52 ±0.09)%	S=1.1	I 67		$h h h' \pi^{\circ} \nu_{\tau}$		(4.76 ±0.06) %	S=1.2
Γ_{15}	$\pi^{-}\pi^{0}$ non- $\rho(770)\nu_{\tau}$		$(3.0 \pm 3.2) \times 10^{-3}$		I 68		$h h h^{\dagger} \pi^{\circ} \nu_{\tau} (\text{ex.} K^{\circ})$		$(4.57 \pm 0.06)\%$	S=1.2
Γ16	$K^{-}\pi^{0}\nu_{\tau}$	[a]	$(4.29 \pm 0.15) \times 10^{-3}$		l ₆₉		$h^{-}h^{-}h^{+}\pi^{0}\nu_{\tau}$ (ex. K ⁰ , ω)		(2.79 ± 0.08)%	S=1.2
	$h^{-} > 2\pi^{0} \mu_{-}$	[-]	$(10.87 \pm 0.11)\%$	S=1.2	Г ₇₀		$\pi^-\pi^+\pi^-\pi^0\nu_{ au}$		(4.62 ± 0.06)%	S=1.2
Г <u>1</u>	$h = 2\pi^0 \nu_{\tau}$		(10.07 ± 0.11) %	S=1.2 S=1.1	Γ ₇₁		$\pi^{-}\pi^{+}\pi^{-}\pi^{0}\nu_{\tau}$ (ex. K^{0})		(4.48 ± 0.06)%	S=1.2
Г 18	$h^{-}2\pi^{0}\mu$ (ex K^{0})		(9.32 ± 0.11) %	S_1.1 S_1.2	Γ ₇₂		$\pi^-\pi^+\pi^-\pi^0 u_ au$ (ex. ${\cal K}^0$, ω)	[a]	(2.70 ± 0.08)%	S=1.2
г 19 Г	$-2 - 2 - 0 + (0 + 10^{\circ})$	[_]	(9.30 ± 0.11) %	5-1.2	Γ ₇₃		$h^- ho \pi^0 u_ au$			
20	$\pi 2\pi^2 \nu_{\tau} (\text{ex.} \text{K}^2)$	a	(9.30 ±0.11)%	5=1.2	Γ ₇₄		$h^- \rho^+ h^- \nu_{\tau}$			
					г ^{, ,}		h = h + h			

For this problem we are interested in:

 $\tau^{\pm} \to \rho^{\pm} \nu \to \pi^{\pm} \pi^0 \nu$ $\tau^{\pm} \to a_1^{\pm} \nu \to 3\pi\nu$

43.8% of the possible decays of the τ fit into these categories.



Józefowicz, Richter-Was, Was https://arxiv.org/abs/1608.02609

- Feature space used:
 - Invariant mass of the intermediate resonances (ρ , a_1).
 - Acoplanarity: angle between the decay planes of the products of each T decay.
 - y (energy asymmetry variables):

$$y_{\rho^{0}}^{\pm} = \frac{E^{\pi^{+}} - E^{\pi^{-}}}{E^{\pi^{+}} + E^{\pi^{-}}}, \quad y_{a_{1}}^{\pm} = \frac{E^{\rho^{0}} - E^{\pi^{\pm}}}{E^{\rho^{0}} + E^{\pi^{\pm}}} - \frac{m_{a_{1}}^{2} - m_{\pi^{\pm}}^{2} + m_{\rho^{0}}^{2}}{2m_{a_{1}}^{2}}$$

• 4-momenta of visible decay products and intermediate resonances: $p^{\mu} = (E, \underline{p})$.

Józefowicz, Richter-Was, Was <u>https://arxiv.org/abs/1608.02609</u>



Contex

22

EXAMPLES: SPIN-PARITY OF THE HIGGS BOSON

Acoplanarity distributions for $H \rightarrow \tau^+ \tau^-$ resulting in 2p's





Józefowicz, Richter-Was, Was https://arxiv.org/abs/1608.02609

- Model: A deep neural network with 6 layers each with 300 perceptrons and a single output perceptron in the 7th layer.
 - Output perceptron uses a sigmoid activation function to compute y_h; all other perceptrons use a ReLU activation function.

Contex

A. Bevan

- Use the Adam optimiser, with batch optimisation and dropout to improve hyperparameter training.
- Use the area under the ROC curve (AUC) as the metric to compare models. Theoretically best result is AUC = 0.782.
- Optimise a loss function that is based on the probability of each of the possible outcomes for each event:

$$-\ln p(y|y_h) = -(y=0)\ln(y_h) - (y=1)\ln(1-y_h)$$

Józefowicz, Richter-Was, Was https://arxiv.org/abs/1608.02609

Aside: ROC: Receiver Operating Characteristic

- A plot of the signal vs noise.
- Integrate the normalised distributions of signal and background predictions as a function of model output; and plot 1-background efficiency vs signal efficiency.
- Illustration here is for 3 different model types; we have encountered the Fisher discriminant. The MLP will be covered later in the course, and the BDT is described elsewhere [1].



[1] For example see <u>https://pprc.qmul.ac.uk/~bevan/teaching/ATLAS-UK-ML.html</u>



Contex

25

EXAMPLES: SPIN-PARITY OF THE HIGGS BOSON

Features/variables	Decay mode: $\rho^{\pm} - \rho^{\mp}$	Decay mode: $a_1^{\pm} - \rho^{\mp}$ $a^{\pm} \rightarrow \rho^0 \pi^{\pm} \rho^0 \rightarrow \pi^+ \pi^-$	Decay mode: $a_1^{\pm} - a_1^{\mp}$ $a^{\pm} \rightarrow a_1^0 \pi^{\pm} - a_1^0 \rightarrow \pi^+ \pi^-$
		$\begin{bmatrix} a_1 \rightarrow \rho \ \pi \ \rho^{\mp} \rightarrow \pi^0 \ \pi^{\mp} \\ & \rho^{\mp} \rightarrow \pi^0 \ \pi^{\mp} \end{bmatrix}$	$\begin{bmatrix} a_1 \rightarrow p \ \pi & p \ \mu & \mu \end{bmatrix} \rightarrow \pi \pi$
$\varphi_{i,k}^*$	1	4	16
$\varphi_{i,k}^*$ and y_i, y_k	3	9	24
$\varphi_{i,k}^*$, 4-vectors	25	36	64
$\varphi_{i,k}^*, y_i, y_k \text{ and } m_i, m_k$	5	13	30
$\varphi_{i,k}^*, y_i, y_k, m_i, m_k$ and 4-vectors	29	45	78

Features/var-	Decay mode: $\rho^{\pm} - \rho^{\mp}$	Decay mode: $a_1^{\pm} - \rho^{\mp}$	Decay mode: $a_1^{\pm} - a_1^{\mp}$
iables	$ ho^{\pm} ightarrow \pi^0 \; \pi^{\pm}$	$a_1^{\pm} \rightarrow ho^0 \pi^{\mp}, \ ho^0 \rightarrow \pi^+ \pi^-$	$\begin{vmatrix} a_1^{\pm} \rightarrow ho^0 \pi^{\pm}, \ ho^0 \rightarrow \pi^+ \pi^- \end{vmatrix}$
		$ ho^{\mp} ightarrow \pi^0 \ \pi^{\mp}$	
True classification	0.782	0.782	0.782
$\varphi_{i,k}^*$	0.500	0.500	0.500
$\varphi_{i,k}^*$ and y_i, y_k	0.624	0.569	0.536
4-vectors	0.638	0.590	0.557
$\varphi_{i,k}^*$, 4-vectors	0.638	0.594	0.573
$\varphi_{i,k}^*, y_i, y_k \text{ and } m_i^2, m_k^2$	0.626	0.578	0.548
$\varphi_{i,k}^*, y_i, y_k, m_i^2, m_k^2$ and 4-vectors	0.639	0.596	0.573

0.5 is a random prediction, 1.0 is a perfect prediction. These models are all somewhere in between, and none match the ideal best possible prediction.

Józefowicz, Richter-Was, Was <u>https://arxiv.org/abs/1608.02609</u>

Conte

26

SUMMARY

- We have discussed the issue of making a decision given the output of a neural network.
 - Binary classification problem using a single perceptron in the output layer.
 - Multiple perceptrons in the output layer for a multi-class output.
 - Use of softmax activation functions as an alternative to multiple perceptrons in the output layer.
- Two examples have been discussed:
 - The MNIST benchmark handwriting recognition problem.
 - \blacktriangleright Determining the J^P of the Higgs boson using $H \to \tau^+ \tau^-$ decays.



SUGGESTED READING

- Discussion of event classification in text books
 - MacKay: Information theory, inference and learning algorithms
 - Chapter: V
 - C. Bishop: Neural Networks for Pattern Recognition
 - Chapter: 1
 - C. Bishop: Pattern Recognition and Machine Learning
 - Chapter: 1
 - > T. Hastie, R. Tibshirani, J. Friedman, Elements of statistical learning
 - Chapter: 2, 4, 9,11
- Examples of classifiers used in particle physics:
 - Bialas, Nemeth, Richter-Wąs, A multi-instance deep neural network classifier: application to Higgs boson CP measurement, <u>arXiv:1803.00838</u>.
 - Madrazo, Cacha, Iglesias, Lucas, Application of a Convolutional Neural Network for image classification to the analysis of collisions in High Energy Physics, <u>arXiv:1708.07034</u>.
 - Abrahão et al., Novel event classification based on spectral analysis of scintillation waveforms in Double Chooz, <u>arXiv:1710.04315</u>.
 - Józefowicz, Richter-Was, Was, Potential for optimizing Higgs boson CP measurement in H → ττ decay at LHC and ML techniques, <u>arxiv:1608.02609</u>.

