

Lecture 3 - Training NN

- Main points; instructor intuition; areas to review.

Recap

- (*) Minimal architecture constraints on NN capabilities
- I/O representation \rightarrow later in course
- Q: How do we compose the network that performs requisite function?

• Perceptron as basic unit; fires/turns output to 1
(threshold)

• Modify from weighted combo \rightarrow weighted combo + bias
affine

BR: what is formal definition for 8th
to be linear (linearity in parameters) : $f(x) = af(x) + bf(y)$
 w

Affine $f(x) : \exists c, g(x) = f(x) - c$
 $g(x)$ - linear

(*) Affine combination of inputs + activation

$$y = \sigma(z) \quad z = \sum_i w_i x_i$$

- Remember modification of bias

- forward architecture: No feedback into inputs.

- Assume network architecture
is sufficient

- Network is a determ. function

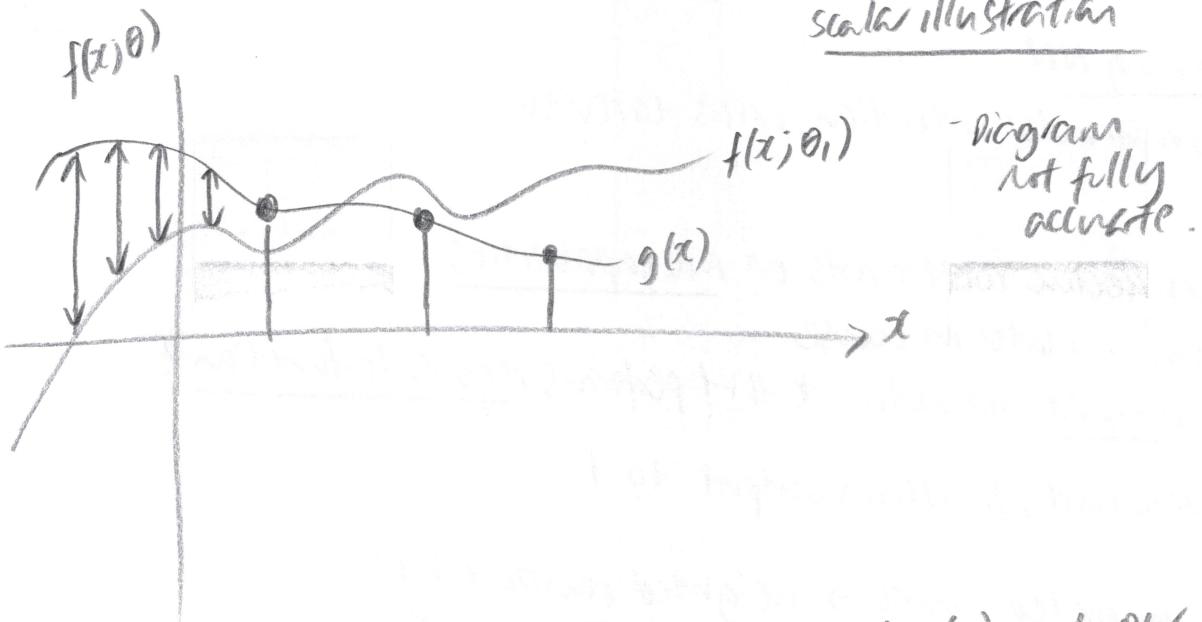
- Parameters of network - weights, biases

- Parameters must be set
(weights, biases)

Q: How to construct MLP to represent a function?

1. By hand (A) \rightarrow check the by-hand intuition

2. Automatic tuning of parameters of network



- BR: quantify gap between function we want to model $g(x)$ and our NN estimated fn ; and area between curves.
- $f(x; w)$ has capacity to rep. $g(x)$
- $\hat{w} = \underset{w}{\operatorname{argmin}} \int_X \operatorname{div}(f(x; w), g(x)) dx$
- This assumes $g(x)$ is known at every x .
 (*) In practice ; $g(x)$ is not specified ; only samples from $g(x)$
- Only noisy observations/samples from $g(x)$ i.e. input-output pairs
- (*) estimate/train NN : determine weights/biases required for it to model a desired function

-
- original MLP (binary classification) (History)
 - fairly comfortable with most of this so far; may not need review
 - restating perceptron
 - Pad data with 1
- $$y = \begin{cases} 1 & \text{if } \sum_{i=1}^{N+1} w_i x_i \geq 0 \\ 0 & \text{otherwise} \end{cases} \quad x_{N+1} = 1$$

- Boundary \rightarrow hyperplane through origin
- can't remember perceptron exactly; but am clear on what happens

A lot of exposition on perceptron

Non-Kor perceptron theorem (bound on convergence) - $(\frac{R}{\gamma})^2$ misclass.

- Bound on classification

Illustration of how perceptron algorithm to model complex decision boundary (Heuristic) (not comput. feasible)

- exponential input search, relabelling

(*) Training MLP using perceptron \rightarrow combinatorial opt.

NP, exponential complexity. , intermediate labels

(*) Do not know outputs of indiv. intermediate neurons in net for train. input

- numerically \rightarrow ADALINE/MADALINE (greedy algorithm)

(W/A1) - Post IEC review - down to

(W/A2) - clarity on exponential time complexity / perceptron failure for complex decision boundary (*)

(*) BR: MLP as combinatorial opt \rightarrow stiffed research

④ Issue lies in threshold function

(*) no way of knowing whether small parameter perturbations reduce error (i.e. an issue to do with information gained from feedback)

- i.e. individual neurons weights can change significantly without changing error.

- to simple MLP: derivatives of function with respect to weights
o nearly everywhere; no derivatives at boundaries

(W/A3) - Review for clarity

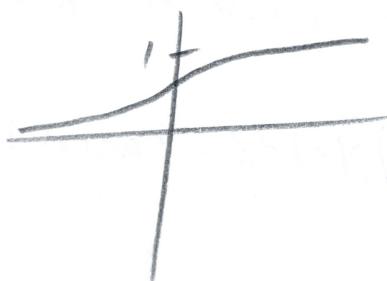
Issue 2: real-life data 'not linearly separable' (~)

Sol: make activation function differentiable

(*) replace threshold activ in perceptron with sigmoid activations:-

- can make $\delta(z)$ approximate threshold

$$\delta(z) = \frac{1}{1+e^{-z}}$$



B

BR: Sigmoid activation

- i) tells me if I'm classifying correctly
- ii) gives an indication of how close I am.

- non-linearly separable data

- even in 1D cases - not possible to threshold counterexamples

- sigmoid probabilistic interpretation $p(y=1|x)$ (recall classification)

- logistic regression: (W) (A) - refresh

↳ Perceptron with sigmoid activation (pr input belongs to class 1)

(R) - key slide

- $\sigma(z)$ - differentiable fn.

$$\frac{dy}{dz} = \sigma'(z)$$

- Diff output y wrt:-

$$\frac{dy}{dw_i} = \frac{dy}{dz} \frac{dz}{dw_i} = \sigma'(z) x_i$$

i) inputs x_i

$$\frac{dy}{dx_i} = \frac{dy}{dz} \frac{dz}{dx_i} = \sigma'(z) w_i$$

ii) param w_i

- via chain rule y is a diff function of both inputs x_i and weights x_i

(*) can compute change in output for small changes in either inputs or weights (W): think tuning knobs physically

some holds for more complex networks.

(W) (S) - review check you understand the examples (quick)

- If activations differentiable, overall function is different. not every param.

BR:

- sets stage for training MLP from input output combos.

$(x_1, d_1); (x_2, d_2) \dots; (x_N, d_N)$

- from earlier; modify

(W) (A) - go over not ~~you~~ some eleven. Bishop

- Assume X is a r.v.

$$\hat{w} = \operatorname{argmin}_w \int_X \operatorname{div}(f(x; w), g(x)) p(x) dx = \operatorname{argmin}_w \underset{x \sim p(x)}{\mathbb{E}} [\operatorname{div}(f(x; w), g(x))]$$

BR: His presentation misses out overfitting, expected error \times

Expected error:

$$\mathbb{E}[\text{div}(f(x; w), g(x))] = \int_X \text{div}(f(x; w), g(x)) p(x) dx$$

BR: Average error
over entire input
space

Empirical est:

$$\mathbb{E}[\text{div}(f(x_i; w), g(x_i))] \approx \frac{1}{N} \sum_{i=1}^N \text{div}(f(x_i; w), d_i)$$

BR: Average error
over samples

- $\textcircled{1}$ $\textcircled{2}$ fit this in with your existing understanding (in Bishop)

(*) minimise empirical risk i.e. go for ERM

given $(x_1, d_1), \dots, (x_N, d_N)$

- error on i^{th} distance: - $\text{div}(f(x_i; w), d_i)$

- empirical average error (empirical risk) on all training data:-

$$\text{loss}(w) = \frac{1}{N} \sum_i \text{div}(f(x_i; w), d_i)$$

- estimate param. to min empirical est of exp. risk:-

$$\hat{w} = \underset{w}{\operatorname{arg\,min}} \text{loss}(w)$$

- min empirical risk over training samples

(non-convex)

(*) instance of function minimisation \rightarrow optimisation

(*) differentiability \rightarrow Δ pacm effect on output \rightarrow effect on error

optimisation crash course

- usual account of a derivative through limits, infinitesimals

- some issues with L'Hopital's rule (he brings up)

- presentation of gradient operator (engineering approach) - I'm okay thanks

$$\nabla_x y = \begin{bmatrix} \frac{\partial y}{\partial x_1} \\ \vdots \\ \frac{\partial y}{\partial x_D} \end{bmatrix}$$

BR

- skips $\textcircled{1}$ $\textcircled{2}$ - review gradients

Q19) - Are you comfortable with rate of increase of function (scalar) and gradient

(*) Line product justif. of gradient

(*) Gradient of a scalar function of a vector points in direction of rate of greatest increase of function; -ve rate of fastest decrease of function