

# Multi-Layer Perceptrons and Deep Learning

Christos Dimitrakakis

November 7, 2023

# Outline

Features and layers

## Algorithms

- Random projection

- Back propagation

- Derivatives

# Layering and features

## Fixed layers

- ▶ Input to layer  $x \in R^n$
- ▶ Output from layer  $\hat{y} \in R^m$ .

## Intermediate layers

Combinations of

- ▶ Linear layer
- ▶ Non-linear activation function.

## Linear layers types

- ▶ Dense
- ▶ Sparse
- ▶ Convolutional

## Activation function

Simple transformations of previous output.

- ▶ Sigmoid
- ▶ Softmax

# Linear layers

Example: Linear regression with  $n$  inputs,  $m$  outputs.

- ▶ Input: Features  $\mathbf{x} \in \mathbb{R}^n$
- ▶ Dense linear layer with  $\boldsymbol{\Theta} \in \mathbb{R}^{m \times n}$
- ▶ Output:  $\hat{\mathbf{y}} \in \mathbb{R}^m$

## Dense linear layer

- ▶ Parameters  $\boldsymbol{\Theta} = \begin{pmatrix} \theta_1 \\ \vdots \\ \theta_m \end{pmatrix}$ ,
- ▶  $\theta_i = [\theta_{i,1}, \dots, \theta_{i,n}]$ ,  $\theta_i$  connects the  $i$ -th output  $y_i$  to the features  $\mathbf{x}$ :

$$y_i = \theta_i \mathbf{x}$$

- ▶ In compact form:

$$\mathbf{y} = \boldsymbol{\Theta} \mathbf{x}$$

# Sigmoid activation

## Example: Logistic regression

- ▶ Input  $\mathbf{x} \in \mathbb{R}^n$
- ▶ Intermediate output:  $z \in \mathbb{R}$ ,

$$z = \sum_{i=1}^n \theta_i x_i.$$

- ▶ Output  $\hat{y} \in [0, 1]$ .

## Definition

This activation ensures we get something we can use as a probability

$$f(z) = 1/[1 + \exp(z)].$$

Now  $P_{\theta}(y = 1|x) = \hat{y}$ .

# Softmax layer

Example: Multivariate logistic regression with  $m$  classes.

- ▶ Input: Features  $\mathbf{x} \in \mathbb{R}^n$
- ▶ Middle: Fully-connected Linear activation layer  $\mathbf{z} = \boldsymbol{\Theta}\mathbf{x}$ .
- ▶ Output:  $\hat{\mathbf{y}} \in \mathbb{R}^m$

## Softmax output layer

We want to translate the real-valued  $z_i$  into probabilities:

$$\hat{y}_i = \frac{\exp(z_i)}{\sum_j \exp(z_j)}.$$

Now  $P_{\boldsymbol{\Theta}}(y = i | \mathbf{x}) = \hat{y}_i$

# Random projections

- ▶ Features  $x$
- ▶ Hidden layer activation  $z$
- ▶ Output  $y$

## Hidden layer: Random projection

Here we project the input into a high-dimensional space

$$z_i = \text{sgn}(\theta_i^\top x) = y_i$$

where  $\Theta = [\theta_i]_{i=1}^m$ .

## The reason for random projections

- ▶ The high dimension makes it easier to learn.
- ▶ The randomness ensures we are not learning something spurious.

# Background on back-propagation

## The problem

- ▶ We need to minimise a loss function  $\ell$
- ▶ We need to calculate

$$\nabla_{\boldsymbol{\theta}} \mathbb{E}_{\boldsymbol{\theta}}[\ell] \approx \frac{1}{T} \sum_{t=1}^T \nabla_{\boldsymbol{\theta}} c(x_t, y_t, \boldsymbol{\theta}).$$

- ▶ However  $c(x_t, y_t, \boldsymbol{\theta})$  is a complex non-linear function of  $\boldsymbol{\theta}$ .

## The solution

- ▶ [1673] Leibniz, the chain rule of differentiation.
- ▶ [1976] Rosenblatt's perceptron without realising it!
- ▶ [1982] Werbos applied it to MLPs.
- ▶ [1986] Rumelhart, Hinton and Williams popularised it.



# Back-propagation

## The chain rule

$$f : X \rightarrow Z, \quad g : Z \rightarrow Y, \quad \frac{dg}{dx} = \frac{dg}{df} \frac{df}{dx}, \quad \nabla_x g = \nabla_f g \nabla_x f$$

## Linear regression

- ▶  $f_{\theta}(x) = \sum_{i=1}^n \theta_i x_i$ .
- ▶  $\mathbb{E}_{\theta}[\ell] \approx \ell(D, \theta) = \frac{1}{T} \sum_{t=1}^T c(\theta, \mathbf{x}_t, y_t)$ .

$$\nabla_{\theta} c(\theta, \mathbf{x}_t, y_t) = \nabla_{\theta} \underbrace{[f_{\theta}(\mathbf{x}_t) - y_t]}_z, \quad g(z) = z^2 \quad (1)$$

$$= \nabla_z g(z) \nabla_f z \nabla_{\theta} f(\mathbf{x}_t) \quad (2)$$

$$= 2[f_{\theta}(\mathbf{x}_t) - y_t] \nabla_f [f_{\theta}(\mathbf{x}_t) - y_t] \nabla_{\theta} f_{\theta}(\mathbf{x}_t) \quad (3)$$

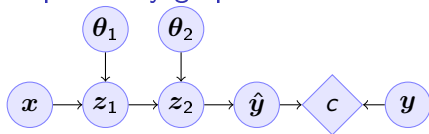
$$= 2[f_{\theta}(\mathbf{x}_t) - y_t] \nabla_{\theta} f_{\theta}(\mathbf{x}_t) \quad (4)$$

# Gradient descent with *back-propagation*

## Inputs

- ▶ Dataset  $D$ , cost function  $\ell = \sum_t c_t$
- ▶ Parametrised architecture with  $k$  layers
  - ▶ Parameters  $\theta_1, \dots, \theta_k$
  - ▶ Intermediate variables:  $z_j = f_j(z_{j-1}, \theta_j)$ ,  $z_0 = x$ ,  $z_k = \hat{y}$ .

## Dependency graph



## Backpropagation with steepest stochastic gradient descent

- ▶ Forward step: For  $j = 1, \dots, k$ , calculate  $z_j = f_j(k)$  and  $c(\hat{y}, y)$
- ▶ Backward step: Calculate  $\nabla_{\hat{y}} c$  and  $d_j \triangleq \nabla_{\theta_j} c = \nabla_{\theta_j} z_j d_{j+1}$  for  $j = k \dots, 1$
- ▶ Apply gradient:  $\theta_j \leftarrow \theta_j - \alpha d_j$ .

# Linear layer

## Definition

This is a linear combination of inputs  $x \in \mathbb{R}^n$  and parameter matrix

$$\Theta \in \mathbb{R}^{m \times n} \text{ where } \Theta = \begin{bmatrix} \theta_1 \\ \vdots \\ \theta_i \\ \vdots \\ \theta_m \end{bmatrix} = \begin{bmatrix} \theta_{1,1} & \cdots & \theta_{1,j} & \cdots & \theta_{1,m} \\ \vdots & \ddots & \vdots & \ddots & \vdots \\ \theta_{i,1} & \cdots & \theta_{i,j} & \cdots & \theta_{i,m} \\ \vdots & \ddots & \vdots & \ddots & \vdots \\ \theta_{n,1} & \cdots & \theta_{n,j} & \cdots & \theta_{n,m} \end{bmatrix}$$

$$f(\Theta, x) = \Theta x \quad f_i(\Theta, x) = \theta_i \cdot x = \sum_{j=1}^n \theta_{i,j} x_j,$$

## Gradient

Each partial derivative is simple:

$$\frac{\partial}{\partial \theta_{i,j}} f_k(\Theta, x) = x_j$$

# Sigmoid layer

$$f(z) = 1/(1 + \exp(-z))$$

## Derivative

So let us ignore the other inputs for simplicity:

$$\frac{d}{dz} f(z) = \exp(-z)/[1 + \exp(-z)]^2$$

# Softmax layer

$$y_i(z) = \frac{\exp(z_i)}{\sum_j \exp(z_j)}$$

## Derivative

$$\frac{\partial}{\partial z_i} y_i(z) = \frac{e^{z_i} e^{\sum_{j \neq i} z_j}}{\left( \sum_j e^{z_j} \right)^2}$$

$$\frac{\partial}{\partial z_i} y_k(z) = \frac{e^{z_i + z_k}}{\left( \sum_j e^{z_j} \right)^2}$$