

Multi-Layer Perceptrons and Deep Learning

Christos Dimitrakakis

October 10, 2024

Outline

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

Linear layers

Example: Linear regression with n inputs, m outputs.

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

Dense linear layer

- ▶ Parameters $\mathbf{B} = \begin{pmatrix} \beta_1 \\ \vdots \\ \beta_m \end{pmatrix}$,
- ▶ $\beta_i = [\beta_{i,1}, \dots, \beta_{i,n}]$, β_i connects the i -th output y_i to the features \mathbf{x} :

$$y_i = \beta_i \mathbf{x}$$
- ▶ In compact form:

$$\mathbf{y} = \mathbf{B}\mathbf{x}$$

Sigmoid activation

Example: Logistic regression

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

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

- ▶ 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 we can interpret $\hat{y} = P_{\beta}(y = 1|\mathbf{x})$.

Softmax layer

Example: Multivariate logistic regression with m classes.

- ▶ Input: Features $\mathbf{x} \in \mathbb{R}^n$
- ▶ Middle: Fully-connected Linear activation layer $\mathbf{z} = \mathbf{B}\mathbf{x}$.
- ▶ 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 we can use $P_B(y = i | \mathbf{x}) = \hat{y}_i$

Random projections

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

Hidden layer: Random projection

Here we project the input into a high-dimensional space

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

where $\mathbf{B} = [\beta_i]_{i=1}^m$, $\beta_{i,j} \sim \text{Normal}(0, 1)$

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 ℓ
- ▶ We need to calculate

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

- ▶ However $c(x_t, y_t, \beta)$ is a complex non-linear function of β .

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

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

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

$$= \nabla_z g(z) \nabla_f z \nabla_{\beta} f(x_t) \quad (2)$$

$$= 2[f_{\beta}(x_t) - y_t] \nabla_f [f_{\beta}(x_t) - y_t] \nabla_{\beta} f_{\beta}(x_t) \quad (3)$$

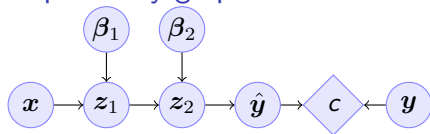
$$= 2[f_{\beta}(x_t) - y_t] \nabla_{\beta} f_{\beta}(x_t) \quad (4)$$

Gradient descent with *back-propagation*

Inputs

- ▶ Dataset D , cost function $\ell = \sum_t c_t$
- ▶ Parametrised architecture with k layers
 - ▶ Parameters β_1, \dots, β_k
 - ▶ Intermediate variables: $z_j = f_j(z_{j-1}, \beta_j)$, $z_0 = x$, $z_k = \hat{y}$.

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 $d_j \triangleq \nabla_{\beta_j} c = \nabla_{\beta_j} z_j d_{j+1}$ for $j = k \dots, 1$
- ▶ Apply gradient: $\beta_j \leftarrow \beta_j - \alpha d_j$.

Other algorithms and gradients

Natural gradient

Defined for probabilistic models

ADAM

Exponential moving average of gradient and square gradients

BFGS: Broyden–Fletcher–Goldfarb–Shanno algorithm

Newton-like method

Example derivatives

Here are some example derivatives

Linear layer

Definition

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

$$\text{where } \mathbf{B} = \begin{bmatrix} \beta_1 \\ \vdots \\ \beta_i \\ \vdots \\ \beta_m \end{bmatrix} = \begin{bmatrix} \beta_{1,1} & \cdots & \beta_{1,j} & \cdots & \beta_{1,m} \\ \vdots & \ddots & \vdots & \ddots & \vdots \\ \beta_{i,1} & \cdots & \beta_{i,j} & \cdots & \beta_{i,m} \\ \vdots & \ddots & \vdots & \ddots & \vdots \\ \beta_{n,1} & \cdots & \beta_{n,j} & \cdots & \beta_{n,m} \end{bmatrix}$$

$$f(\mathbf{B}, \mathbf{x}) = \mathbf{B}\mathbf{x} \quad f_i(\mathbf{B}, \mathbf{x}) = \beta_i \cdot \mathbf{x} = \sum_{j=1}^n \beta_{i,j} x_j,$$

Gradient

Each partial derivative is simple:

$$\frac{\partial}{\partial \beta_{i,j}} f_k(\mathbf{B}, \mathbf{x}) = x_j$$

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

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

Derivative

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

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