

Chapter 2 – Neural Networks

Fabien Petit

Department of Economics, University of Barcelona
Centre for Education Policy and Equalising Opportunities, UCL

Introduction

Neural networks are a family of flexible predictive statistical models with a network structure

- ▶ This lecture introduces the most fundamental model structure, called the **feed-forward network**
- ▶ Specialized architectures for image (CNN), language (RNN), generative (GAN) modelling, Transformers, etc.
- ▶ State-of-the-art performance in problems involving high-dimensional and unstructured data

Introduction

Networks are called “neural” because their basic structure and functioning were loosely inspired by neuroscience

- ▶ This analogy is no longer relevant as developments follow mathematical and engineering principles
- ▶ Early implementations in the 1980's but networks have become popular over the past decade
- ▶ Computing power and algorithmic innovations made networks more capable and easier to estimate

Structure

Units

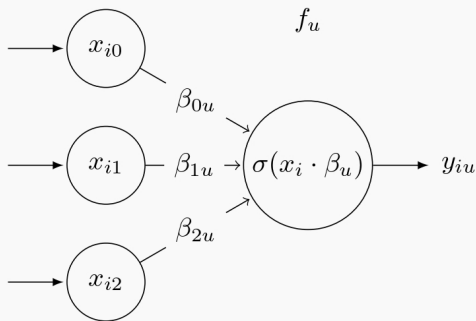
A network is composed of many **simple non-linear functions** called **neurons** or **units**, denoted by u

$$y_{iu} = \sigma(x_{i \leftarrow x_0} \cdot \beta_u) \quad (1)$$

where $x_{i \leftarrow x_0}$ and β_u are vectors of length $k+1$, \leftarrow denotes concatenation and σ is a non-linear function called **activation**

- ▶ A unit performs a dot product between a vector of **variables** and a vector of **parameters**
- ▶ Early networks used a **sigmoid activation function**, so a unit can be seen as a logistic regression model

Units



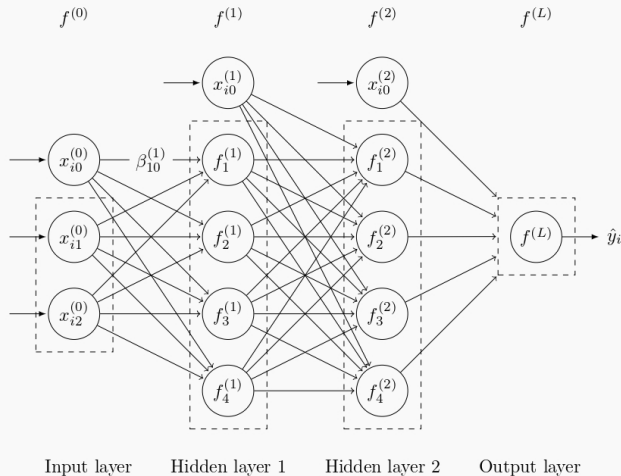
Networks can be represented as directed graphs. Circles are network units, simple arrows are parameters, and plain arrows are the unit's inputs and outputs. For reasons that will become clear, input variables are represented as units despite the absence of computation, i.e., a unit is a number.

Layers

The unit's output is not the model's final prediction but an **intermediate transformation** of the input

- ▶ Multiple transformations are computed using nested groups of units called **layers** $f^{(l)}, l = 0, \dots, L$
- ▶ **Hidden layers** $f^{(1)}, \dots, f^{(L-1)}$ compute recursively intermediate transformations of the input
- ▶ The **output layer** $f^{(L)}$ produces the model prediction in the transformed inputs i.e., the output of $f^{(L-1)}$

Structure



Feed-forward network with three layers

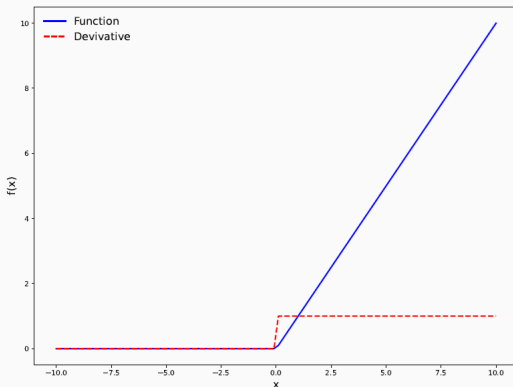
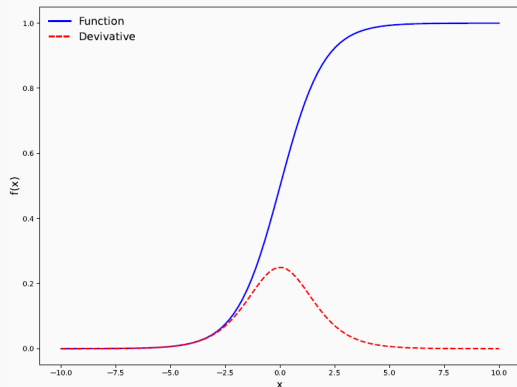
Dashed rectangles represent network layers. The model has two hidden layers with $k^{(1)} = k^{(2)} = 4$ hidden units. Each of the $k^{(l)}$ units of a layer transforms the same inputs using different sets of parameters.

Hidden Layers

Hidden units produce **intermediate transformations** by computing non-linear combinations of their input

- ▶ This allows the optimized model to detect non-linearities and interactions (i.e., functional form)
- ▶ The structure of the hidden layers (i.e., the number of layers and units) is a **tuning parameter**
- ▶ The activation function is chosen for **optimization** (e.g., differentiable, monotonous, zero-centered)

Hidden Layers



Sigmoid (left) and Rectified Linear Units (right). The plain lines represent the activation functions, and their derivatives are depicted by the dashed lines. With linear activation, the network would simplify to a linear model. Many differentiable functions can be used in the hidden layers, but some have computational advantages, e.g., differentiable, monotonous, and zero-centered.

Output Layer

The **output unit** transforms the output of the last hidden layer (i.e., transformed variables) into the **estimated response**

- ▶ The number of output units and their activation depend on the distribution of the response
- ▶ The activations are often the same as the GLM's link functions (e.g., identity, logistic, Poisson)
- ▶ The output layer can be interpreted as a generalized linear model on the transformed variables

Learning

Why would we expect this layered model structure to increase predictive performance?

- ▶ The composition of numerous parametrized non-linear functions gives the model considerable flexibility
- ▶ Unlike linear models, networks do not attempt to fit the entire data space in a single parametric equation
- ▶ They gradually approach a functional form by composing numerous simple non-linear transformations

Learning

Researchers have stressed the importance of distributed and hierarchical representations (Bengio et al. 2013)

- ▶ Network representations or transformed variables can be seen as both **composite** and **latent** variables
- ▶ A hidden unit may combine surface area and the number of bedrooms to determine the type of housing
- ▶ Another may combine the distance to school and the student-teacher ratio into measures of schooling quality
- ▶ The same applies to any feature present in the data, which the model identifies as relevant for prediction

Learning

The units in the second hidden layer may combine these transformed inputs into **more abstract representations**

- ▶ E.g., composite measures of house characteristics, local labor markets, local amenities, or accessibility
- ▶ Hidden units are sub-models constructing increasingly abstract representations from the input
- ▶ Units encode multiple such concepts in a single numerical value, so they are **not readily interpretable**

Generalization

Generalization

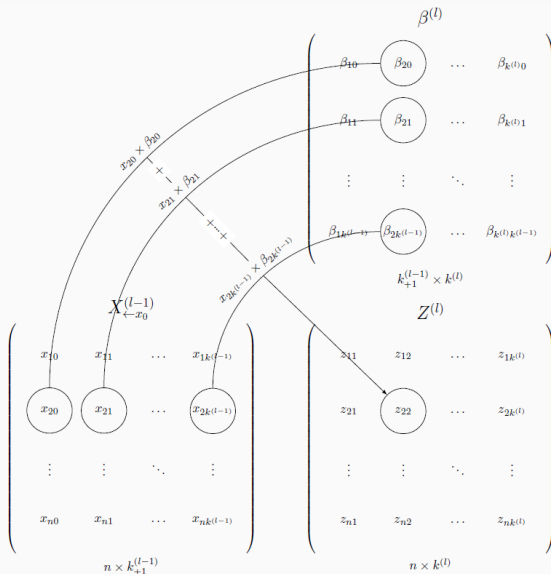
Equation (1) can be generalized to multiple **observations** and **units**. Note that $X^{(l)} = y^{(l)}$ to simplify notation

$$X^{(l)} = y^{(l)} = \sigma^{(l)} \left(X_{\leftarrow x_0}^{(l-1)} \cdot \beta^{(l)} \right), \quad (2)$$

where $X^{(l)}$ and $\beta^{(l)}$ are matrices of concatenated observations (as rows) and parameters (as columns)

$$X^{(l)} = \begin{bmatrix} - & x_1^{(l)} & - \\ & \vdots & \\ - & x_n^{(l)} & - \end{bmatrix}_{n \times k^{(l)}} \quad \beta^{(l)} = \begin{bmatrix} | & & | \\ \beta_0^{(l)} & \dots & \beta_{k^{(l)}}^{(l)} \\ | & & | \end{bmatrix}_{k_{+1}^{(l-1)} \times k^{(l)}}$$

Generalization



Matrix Product

$X_{\leftarrow x_0}^{(l-1)}$ contains n observations (rows) and $k_{+1}^{(l-1)}$ variables (columns), including the constant. $\beta^{(l)}$ contains the corresponding $k_{+1}^{(l-1)}$ parameters for each of the $k^{(l)}$ units (columns). The product of these two matrices is denoted $Z^{(l)}$ with n observations, where each $k^{(l)}$ variable represents a combination of the same input with different parameters. The activation function $\sigma^{(l)}$ is applied element-wise to produce the matrix $X^{(l)}$.

Generalization – Forward Propagation

Operations			Dimensions
$f^{(1)}$	$X^{(1)} = \sigma^{(1)} \left(X_{\leftarrow x_0}^{(0)} \cdot \beta^{(1)} \right)$		$n \times 4 = [n \times 2_{+1}] \cdot [3 \times 4]$
$f^{(2)}$	$X^{(2)} = \sigma^{(2)} \left(X_{\leftarrow x_0}^{(1)} \cdot \beta^{(2)} \right)$		$n \times 4 = [n \times 4_{+1}] \cdot [5 \times 4]$
$f^{(L)} \quad \hat{y} =$	$X^{(L)} = \sigma^{(L)} \left(X_{\leftarrow x_0}^{(2)} \cdot \beta^{(L)} \right)$		$n \times 1 = [n \times 4_{+1}] \cdot [5 \times 1]$

This table details the forward propagation operations for the network illustrated in Figure 2, along with the matrix dimensions. More generally, $X^{(l)}$ has dimensions $n \times k^{(l)}$ while $\beta^{(l)}$ has dimensions $k_{+1}^{(l-1)} \times k^{(l)}$.

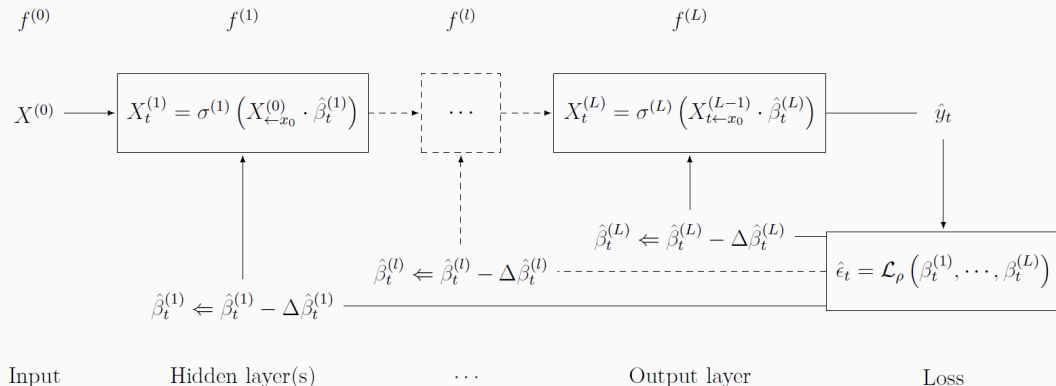
Estimation

Intuition

Since hidden units have **no target output**, the estimated error is evaluated **at the end of the forward pass**

- ▶ The training error is computed using the **loss function** corresponding to the response distribution
- ▶ **Regularization** ρ may sum the parameters, either squared (i.e., Ridge L_2) or expressed as an absolute value (i.e., Lasso L_1)
- ▶ There are many other ways to regularize the model (e.g., early stopping, dropout, normalization, adversarial)

Abstract Feed-Forward Network



The subscript t indicates the quantities that change with each iteration of the optimization routine, \Leftarrow is the update operator and $\Delta \hat{\beta}_t^{(l)}$ is the amount by which the parameters are updated. The loss is computed at the end of the forward pass. The backward pass updates the parameters iteratively to decrease the training error.

Intuition

Within this parameter space, the **optimization** searches the $\hat{\beta}$ that **minimize the (penalized) training error**

- ▶ The loss function is never strictly convex, and the optimization problem has **no closed-form solution**
- ▶ **Gradient-based optimization** updates the parameters iteratively to **converge toward a local minimum**
- ▶ Gradients are computed for each training observation and averaged before updating the parameters

Update

For simplicity, consider a single training observation i and a single parameter of unit u indexed j . The update rule is

$$\beta_j \Leftarrow \beta_j - \eta \frac{\partial \mathcal{L}_i}{\partial \beta_j} \quad (3)$$

where \Leftarrow is the update operator and η is a **tuning parameter** called the **learning rate**

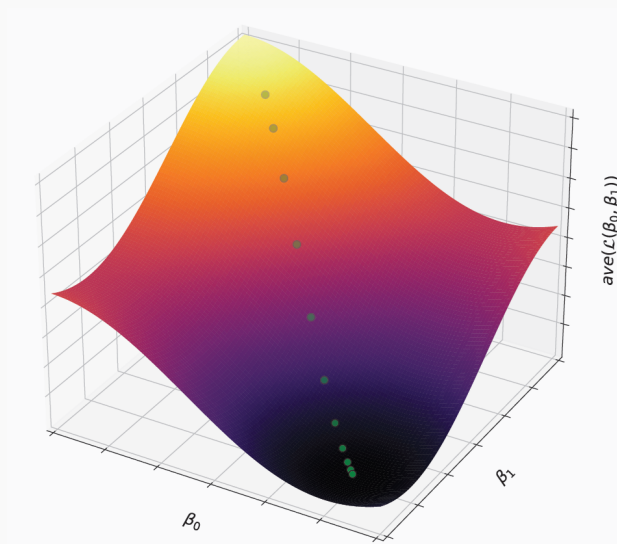
- ▶ Parameters are updated according to their relative **contribution to the training error**
- ▶ This is measured by the partial derivatives of the loss function with respect to the parameters (i.e., gradient)

Optimization Path

The closest local approximation of a function is the **tangent**, whose slope is the partial derivative

- ▶ The vector of partial derivatives with respect to every function parameter is called **the gradient**
- ▶ The gradient points to the **direction** where the function is increasing by the largest amount
- ▶ The **update** is **proportional to the partial derivative** but smaller in magnitude (i.e., local approximation)

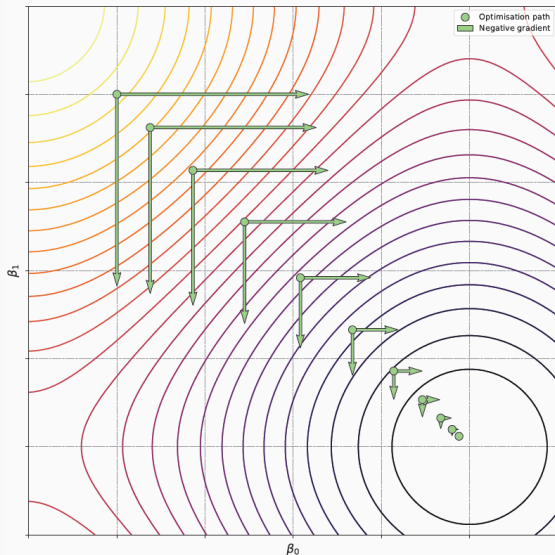
Optimization Path



Loss function with two parameters

Illustrative function with two parameters, but the intuition applies to higher dimensions. The loss function averaged across training observations can be represented as a hilly landscape in the multi-dimensional space of parameter values. Note the presence of local minima and saddle points.

Optimization Path



Loss function with two parameters

In the vector space, the gradient points to the direction where the loss function increases by the largest amount. The steps represent the negative gradient, weighted by the learning rate. The color denotes the loss averaged over the training batch.

Generalization

Generalizing to **multiple training observations**, n partial derivatives are computed for each parameter

- ▶ A change in a single parameter affects the training error for every observation, either positively or negatively
- ▶ Partial derivatives are averaged out to compute the update that reduces the average training error
- ▶ The size of the training data and the number of parameters affect the computational cost of the update

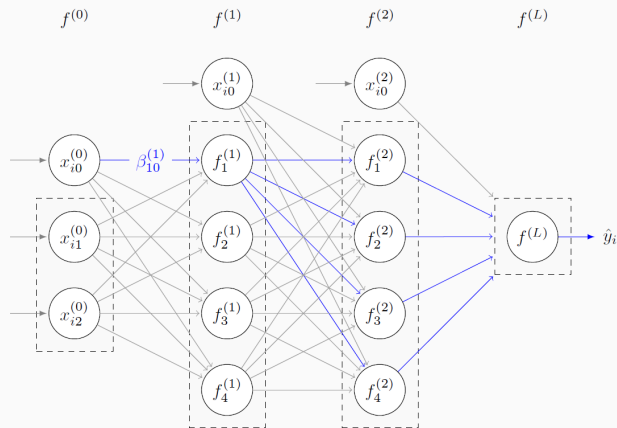
Backpropagation

Backpropagation

Computing the gradients is demanding, given the number of parameters and the nested nature of the model

- ▶ Backpropagation (Rumelhart et al. 1986) addresses these issues using the **chain rule of differentiation**
- ▶ The complex gradient expressions can be expressed as the **product of simple partial derivatives**
- ▶ Duplicate calculations are avoided by using quantities that are computed during the forward pass

Intuition



Backpropagation

A small change to $\beta_{10}^{(1)}$ alters the result of the dot product $z_1^{(1)}$ and the unit's output $x_1^{(1)}$. This modifies the output of the four units in the second hidden layer, the predicted response and training error. To compute the contribution of a parameter, we take into account its impact on all units in the next layers.

Chain Rule

Consider the function composition

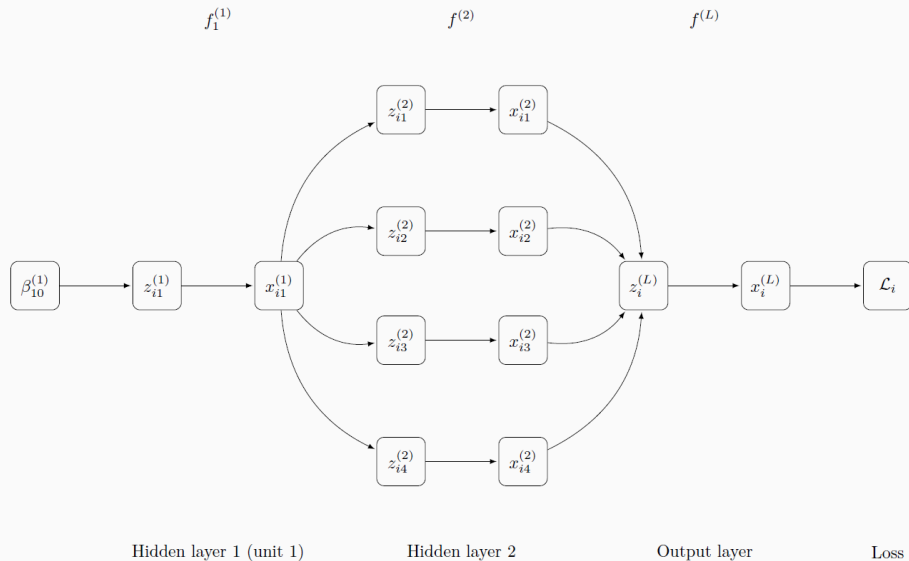
$$\begin{aligned}y &= f(x) \\ z &= g(y) = (g \circ f)(x)\end{aligned}$$

The chain rule states that

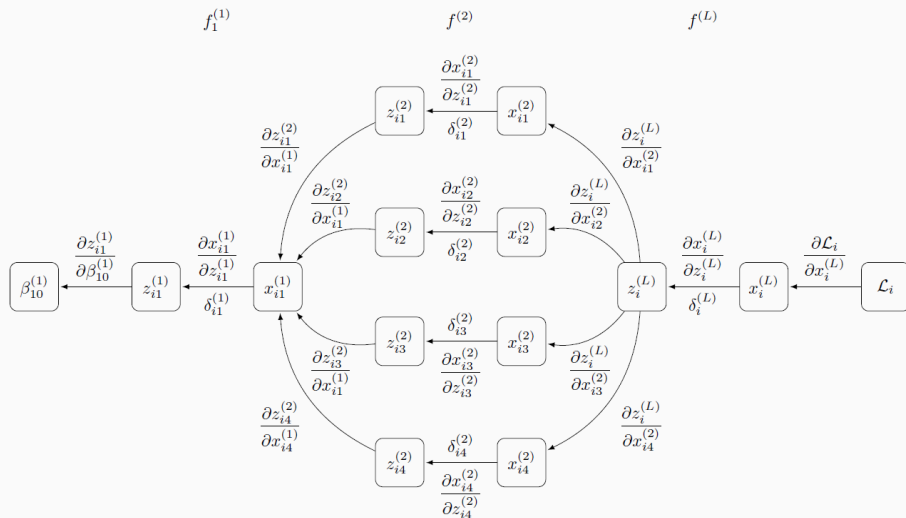
$$\frac{\partial z}{\partial x} = \frac{\partial z}{\partial y} \frac{\partial y}{\partial x}$$

The complex gradient expressions can be expressed as the product of simple partial derivatives

Chain rule – Impact of Changing $\beta_{10}^{(1)}$ on Network



Chain rule – Detailed Backpropagation Operations



Hidden layer 1 (unit 1)

Hidden layer 2

Output layer

Loss

Chain Rule

The partial derivative of the loss with respect to $\beta_{10}^{(1)}$ for observation i can be expressed as

$$\begin{aligned} \frac{\partial \mathcal{L}_i}{\partial \beta_{10}^{(1)}} &= \underbrace{\frac{\partial \mathcal{L}_i}{\partial \mathbf{x}_i^{(L)}} \frac{\partial \mathbf{x}_i^{(L)}}{\partial \mathbf{z}_i^{(L)}}}_{\delta_i^{(L)}} \\ &\quad \underbrace{\frac{\partial \mathbf{z}_i^{(L)}}{\partial \mathbf{x}_{i1}^{(2)}} \frac{\partial \mathbf{x}_{i1}^{(2)}}{\partial \mathbf{z}_{i1}^{(2)}}}_{\prod \dots = \delta_{i1}^{(2)}} \underbrace{\frac{\partial \mathbf{z}_i^{(L)}}{\partial \mathbf{x}_{i2}^{(2)}} \frac{\partial \mathbf{x}_{i2}^{(2)}}{\partial \mathbf{z}_{i2}^{(2)}}}_{\prod \dots = \delta_{i2}^{(2)}} \underbrace{\frac{\partial \mathbf{z}_i^{(L)}}{\partial \mathbf{x}_{i3}^{(2)}} \frac{\partial \mathbf{x}_{i3}^{(2)}}{\partial \mathbf{z}_{i3}^{(2)}}}_{\prod \dots = \delta_{i3}^{(2)}} \underbrace{\frac{\partial \mathbf{z}_i^{(L)}}{\partial \mathbf{x}_{i4}^{(2)}} \frac{\partial \mathbf{x}_{i4}^{(2)}}{\partial \mathbf{z}_{i4}^{(2)}}}_{\prod \dots = \delta_{i4}^{(2)}} \\ &\quad \underbrace{\frac{\partial \mathbf{z}_{i1}^{(2)}}{\partial \mathbf{x}_{i1}^{(1)}} \frac{\partial \mathbf{z}_{i2}^{(2)}}{\partial \mathbf{x}_{i1}^{(1)}} \frac{\partial \mathbf{z}_{i3}^{(2)}}{\partial \mathbf{x}_{i1}^{(1)}} \frac{\partial \mathbf{z}_{i4}^{(2)}}{\partial \mathbf{x}_{i1}^{(1)}} \frac{\partial \mathbf{x}_{i1}^{(1)}}{\partial \mathbf{z}_{i1}^{(1)}} \frac{\partial \mathbf{z}_{i1}^{(1)}}{\partial \beta_{10}^{(1)}}}_{\dots \delta_{i1}^{(1)}} \end{aligned}$$

where $\prod \dots$ denotes the product of the bracket with the lines above and $\delta_{iu}^{(l)}$ is the “error signal” defined as $\partial \mathcal{L}_i / \partial \mathbf{z}_{iu}^{(l)}$

Generalization

More generally, the partial derivative of the loss with respect to $\beta_{ju}^{(l)}$ for observation i can be expressed as

$$\frac{\partial \mathcal{L}_i}{\partial \beta_{ju}^{(l)}} = \frac{\partial \mathcal{L}_i}{\partial z_{iu}^{(l)}} \frac{\partial z_{iu}^{(l)}}{\partial \beta_{ju}^{(l)}} = \delta_{iu}^{(l)} \frac{\partial z_{iu}^{(l)}}{\partial \beta_{ju}^{(l)}} \quad (5)$$

where the “error signal” $\delta_{iu}^{(l)}$ measures the contribution of $z_{iu}^{(l)}$ to the training error

- ▶ Backpropagation computes efficiently the error signals $\delta_{iu}^{(l)}$ for each unit and observation
- ▶ The error signals can be easily related to the partial derivative with respect to each parameter

Generalization

To compute the error signals efficiently, **backpropagation avoids duplicated computations**

- ▶ The error signal of a layer l can be expressed as a function of the error signal in the next layer $l + 1$
- ▶ These error signals are computed starting from the output layer all the way to the first hidden layer
- ▶ For increased efficiency, the relevant calculations can be stored in memory during the forward pass

Generalization – Forward and Backward Passes

Table: Forward and Backward Passes

	Forward propagation		Backpropagation	
$f^{(1)}$	\downarrow	$X^{(0)} = X_{\leftarrow x_0}^{(0)}$	\leftarrow	
	\downarrow	$Z^{(1)} = X^{(0)} \cdot \beta^{(1)}$	\uparrow	$\partial Z^{(1)} / \partial \beta^{(1)} = X^{(0)T}$
	\downarrow	$X^{(1)} = \sigma(Z^{(1)})$	$\delta^{(1)}$	$\partial X^{(1)} / \partial Z^{(1)} = \sigma'(Z^{(1)})$
	\downarrow	$X^{(1)} = X_{\leftarrow x_0}^{(1)}$	\uparrow	$\partial Z^{(2)} / \partial X^{(1)} = \beta^{(2)T}$
\vdots			\vdots	\vdots
$f^{(l)}$	\downarrow	$Z^{(l)} = X^{(l-1)} \cdot \beta^{(l)}$	\uparrow	$\partial Z^{(l)} / \partial \beta^{(l)} = X^{(l-1)T}$
	\downarrow	$X^{(l)} = \sigma(Z^{(l)})$	$\delta^{(l)}$	$\partial X^{(l)} / \partial Z^{(l)} = \sigma'(Z^{(l)})$
	\downarrow	$X^{(l)} = X_{\leftarrow x_0}^{(l)}$	\uparrow	$\partial Z^{(l+1)} / \partial X^{(l)} = \beta^{(l+1)T}$
\vdots			\vdots	\vdots
$f^{(L)}$	\downarrow	$Z^{(L)} = X^{(L-1)} \cdot \beta^{(L)}$	\uparrow	$\partial Z^{(L)} / \partial \beta^{(L)} = X^{(L-1)T}$
	\downarrow	$X^{(L)} = \sigma(Z^{(L)})$	$\delta^{(L)}$	$\partial X^{(L)} / \partial Z^{(L)} = \sigma'(Z^{(L)})$
	\rightarrow		\uparrow	$\partial \mathcal{L}(\beta) / \partial X^{(L)} = \mathcal{L}'(X^{(L)})$

Most of the quantities required to compute the gradients, such as $X^{(l)}$, $\beta^{(l)}$ and $Z^{(l)}$ are stored during forward propagation to avoid duplicated computations

Generalization

The matrix $\delta^{(l)}$ has dimensions $n \times k^{(l)}$ and contains the error signals for every unit and observations

$$\begin{aligned}\delta^{(L)} &= \frac{\partial \mathcal{L}(\beta)}{\partial X^{(L)}} \frac{\partial X^{(L)}}{\partial Z^{(L)}} \\ &= \mathcal{L}' \left(X^{(L)} \right) \odot \sigma^{(L)} \left(Z^{(L)} \right)\end{aligned}\tag{6}$$

$$\begin{aligned}\delta^{(l)} &= \delta^{(l+1)} \frac{\partial Z^{(l+1)}}{\partial X^{(l)}} \frac{\partial X^{(l)}}{\partial Z^{(l)}} \\ &= \delta^{(l+1)} \cdot \beta^{(l+1)T} \odot \sigma'^{(l)} \left(Z_{\leftarrow}^{(l)} x_0 \right)\end{aligned}\tag{7}$$

where \odot denotes the Hadamard product and the gradient is represented as a column vector (i.e., denominator layout)

Generalization

Using the expression for the error signal, the matrix of partial derivatives with respect to each parameter is computed as

$$\frac{\partial \mathcal{L}(\beta)}{\partial \beta^{(l)}} = X_{\leftarrow x_0}^{(l-1)T} \cdot \delta_{\rightarrow x_0}^{(l)} \quad (8)$$

For each parameter, this formulation sums the gradients across training observations. The update rule becomes

$$\beta^{(l)} \Leftarrow \beta^{(l)} - \frac{\eta}{n} \frac{\partial \mathcal{L}(\beta)}{\partial \beta^{(l)}} \quad (9)$$

In practice, each gradient descent iteration is performed using a partition of the training observations called a “**batch**”

Generalization

	Operations	Dimensions
$f^{(L)}$	$\delta^{(L)} = \mathcal{L}'(X^{(L)}) \odot \sigma^{(L)'}(Z^{(L)})$ $\frac{\partial \mathcal{L}(\beta)}{\partial \beta^{(L)}} = X_{\leftarrow x_0}^{(2)T} \cdot \delta^{(L)}$	$n \times 1 = [n \times 1] \odot [n \times 1]$ $5 \times 1 = [4_{+1} \times n] \cdot [n \times 1]$
$f^{(2)}$	$\delta^{(2)} = \delta^{(L)} \cdot \beta^{(L)T} \odot \sigma^{(2)'}(Z_{\leftarrow x_0}^{(2)})$ $\frac{\partial \mathcal{L}(\beta)}{\partial \beta^{(2)}} = X_{\leftarrow x_0}^{(1)T} \cdot \delta_{\rightarrow x_0}^{(2)}$	$n \times 5 = [n \times 1] \cdot [1 \times 5] \odot [n \times 4_{+1}]$ $5 \times 4 = [4_{+1} \times n] \cdot [n \times 5_{-1}]$
$f^{(1)}$	$\delta^{(1)} = \delta_{\rightarrow x_0}^{(2)} \cdot \beta^{(2)T} \odot \sigma^{(1)'}(Z_{\leftarrow x_0}^{(1)})$ $\frac{\partial \mathcal{L}(\beta)}{\partial \beta^{(1)}} = X_{\leftarrow x_0}^{(0)T} \cdot \delta_{\rightarrow x_0}^{(1)}$	$n \times 5 = [n \times 5_{-1}] \cdot [4 \times 5] \odot [n \times 4_{+1}]$ $3 \times 4 = [2_{+1} \times n] \odot [n \times 5_{-1}]$

The denominator layout implies the inversion of the two terms and the transpose for the matrix multiplications. Note that $\delta^{(l+1)} \cdot \beta^{(l+1)T}$ contains the bias while $\sigma' \left(Z^{(l)} \right)$ does not. To allow for element-wise multiplication, a column of ones is added to $Z^{(l)}$. When we back-propagate to the previous layer, $\delta^{(l)}$ contains the error linked to the bias, which is removed because it is not connected to the previous layer.

Gradient Checking

You can check your implementation of backpropagation against a numerical approximation of the gradient

$$\frac{\partial \mathcal{L}(\beta)}{\partial \beta} \approx \frac{\mathcal{L}(\beta + \epsilon) - \mathcal{L}(\beta - \epsilon)}{2\epsilon}$$

where ϵ is a small random value. In the case where β is a vector of values, we can compute for each β_j , $j = 1, \dots, k$

$$\frac{\partial \mathcal{L}(\beta_3)}{\partial \beta_0} \approx \frac{\mathcal{L}(\beta_0 + \epsilon, \beta_1, \dots, \beta_k) - \mathcal{L}(\beta_0 - \epsilon, \beta_1, \dots, \beta_k)}{2\epsilon}$$

$$\frac{\partial \mathcal{L}(\beta_1)}{\partial \beta_1} \approx \frac{\mathcal{L}(\beta_0, \beta_1 + \epsilon, \dots, \beta_k) - \mathcal{L}(\beta_0, \beta_1 - \epsilon, \dots, \beta_k)}{2\epsilon}$$

...

Application

Hedonic Pricing

This network can be used as a drop-in replacement for linear models **when interpretability is not required**

- ▶ This is illustrated on a hedonic price problem on the Boston dataset (Harrison and Rubinfeld 1978)
- ▶ The response is the median value of owner-occupied houses for 506 census tracts in Boston in 1970
- ▶ The input contains 13 variables, including measures of air quality, schooling quality, among others

Hedonic Pricing – Comparison of Predictive Models

Model	R^2 training		R^2 test	
Log-linear regression	0.7416	(0.0091)	0.7144	(0.0743)
Feed-forward network	0.9908	(0.0015)	0.8664	(0.0434)
- regularised	0.9502	(0.0023)	0.9016	(0.0498)

Statistics are produced using 10-fold cross-validation. Standard deviations are in parentheses. The network has 2 hidden layers with 64 units and ReLU activation. The output layer has a single unit with sigmoid activation. Note that the (irreducible) error is substantial in social science data.

Summary

Summary

The processing of **high-dimensional and unstructured data** offers strong potential for original economic analysis

- ▶ Traditional models perform poorly due to the dimensionality and the absence of theory
- ▶ Neural networks are **powerful and flexible predictive models** designed to solve these problems
- ▶ Specialized networks offer state-of-the-art performance on image and language data, among other sources

Summary

On the downside, networks **trade most interpretability for flexibility** and are **computationally heavier**

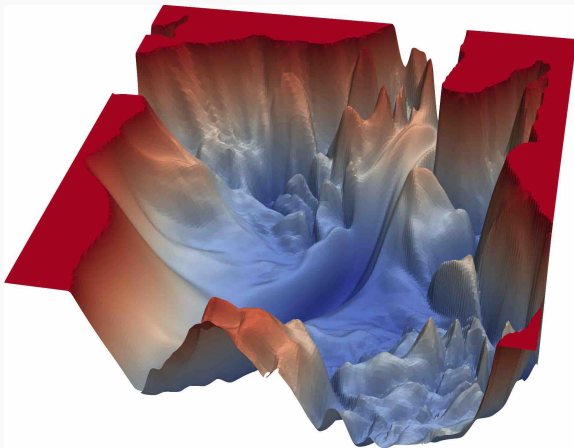
- ▶ Approximating even simple functions involves many parameters and requires more observations
- ▶ Many **tuning parameters**, e.g., layers, units, activation, loss, optimizer, regularization, learning rate
- ▶ Network structures and implementation details are important in this more empirical literature

Summary

Computing the gradients is demanding, given the number of parameters and the nested nature of the model

- ▶ The **backpropagation** algorithm addresses these issues by using the **chain rule of differentiation**
- ▶ The complex gradient expressions can be expressed as a product of simple partial derivatives
- ▶ It avoids duplicate calculations by reusing quantities that are computed and stored during the forward pass

Summary – Illustrative Loss Landscape



Illustrative loss landscape

Three-dimensional representation of a VGG56 network loss landscape (Li et al. 2018). Complex loss landscapes have many local minima and saddle points where the optimization could get trapped and converge to sub-optimal values

References

- ▶ Harrison, David J. and Daniel L. Rubinfeld (1978). "Hedonic housing prices and the demand for clean air". *Journal of Environmental Economics and Management* 5.1, pp. 81-102 (cit. on p. 44).
- ▶ Rumelhart, David E., Geoffrey E. Hinton, and Ronald J. Williams (1986). "Learning representations by back-propagating errors". *Nature* 323.6088, pp. 533-536 (cit. on p. 30).
- ▶ Bengio, Yoshua, Aaron Courville, and Pascal Vincent (2013). "Representation learning: A review and new perspectives". *IEEE Transactions on Pattern Analysis and Machine Intelligence* 35.8, pp. 1798-1828 (cit. on p. 15).
- ▶ Nielsen, Michael A. (2015). *Neural networks and deep learning*. Determination Press.
- ▶ Goodfellow, Ian, Yoshua Bengio, and Aaron Courville (2016). *Deep Learning*. MIT Press.
- ▶ Li, Hao et al. (2018). "Visualizing the loss landscape of neural nets". *Advances in Neural Information Processing Systems*. Vol. 31 (cit. on p. 50).