

Contents

Introduction	5
1 Presentation of the field	7
1.1 Tensors	8
1.1.1 Definition	8
1.1.2 Manipulation	9
1.1.3 Binary operations	12
1.2 Neural Networks	15
1.2.1 Simple formalization	15
1.2.2 Generic formalization	18
1.2.3 Interpretation	20
1.2.4 Parameterization and training	21
1.2.5 Examples of layer	24
1.2.6 Examples of regularization	29
1.2.7 Examples of architecture	29
1.3 Graphs and signals	31
1.3.1 Basic definitions	31
1.3.1.1 Graphs	31
1.3.1.2 (Real-valued) Signals	33
1.3.2 Graphs in deep learning	34
1.3.2.1 Connectivity graph	35
1.3.2.2 Computation graph	37

1.3.2.3	Underlying graph structure and signals	37
1.3.2.4	Graph-structured dataset	37
2	Convolution on graph domains	39
2.1	Analysis of the classical convolution	40
2.1.1	Properties of the convolution	40
2.1.2	Characterization on grid graphs	40
2.1.3	Usefulness of convolutions in deep learning	43
2.1.4	Related work generalising the discrete convolution on graphs	44
2.2	Construction on the vertex set through groups	45
2.3	Construction with the edge set	53
2.3.1	Construction on Cayley graphs	53
2.3.2	Construction on graph groupoids	54
2.3.3	To rename	57
2.3.4	Lattice-regular graph	58
2.4	Conclusion of chapter 2	59
3	Neural networks on graphs	61
4	Industrial applications	63
	Conclusion	65
	Bibliography	67
0	Trash bin and more drafts	75
0.1	Naming conventions	76
0.1.1	Basic notions	76
0.1.2	Graphs and graph signals	76
0.1.3	Data and datasets	77
0.2	Disambiguation of the subject	77

<i>CONTENTS</i>	3
0.2.1 Irregularly structured data	78
0.2.2 Unstructured data	79
0.3 Datasets	79
0.4 Tasks	79
0.5 Goals	79
0.6 Invariance	79
0.7 Methods	79
0.8 Expressivity analysis of dense versus sparse connectivity . . .	80
0.8.1 Strong regular case	80
0.8.2 Draft	81
0.9 Conv drafts	82
0.9.1	85
0 Tentative previsional plans	86
0 Keywords and tentative titles	93

Introduction

TODO:

Chapter 1

Presentation of the field

In this section, we present notions related to our domains of interest. In particular, for tensors we give original definitions that are more appropriate for our study. In the neural network's section, we present the concepts necessary to understand the evolution of the state of the art research in this field. In the last section, we present graphs for their usage in deep learning.

Vector spaces considered in what follows are assumed to be finite-dimensional and over the field of real numbers \mathbb{R} .

1.1	Tensors	8
1.2	Neural Networks	15
1.3	Graphs and signals	31

1.1 Tensors

Intuitively, tensors in the field of deep learning are defined as a generalization of vectors and matrices, as if vectors were tensors of rank 1 and matrices were tensors of rank 2. That is, they are objects in a vector space and their dimensions are indexed using as many indices as their rank, so that they can be represented by multidimensional arrays. In mathematics, a tensor can be defined as a special type of multilinear function (Bass, 1968; Marcus, 1975; Williamson, 2015), which image of a basis can be represented by a multidimensional array. Alternatively, Hackbush propose a mathematical construction of a tensor space as a quotient set of the span of an appropriately defined tensor product (Hackbusch, 2012), which coordinates in a basis can also be represented by a multidimensional array. In particular in the field of mathematics, tensors enjoy an intrinsic definition that neither depend on a representation nor would change the underlying object after a change of basis, whereas in our domain, tensors are confounded with their representation.

1.1.1 Definition

Our definition of tensors is such that they are a bit more than multidimensional arrays but not as much as mathematical tensors, for that they are embedded in a vector space so that deep learning objects can be later defined rigorously.

Given canonical bases, we first define a tensor space, then we relate it to the definition of the tensor product of vector spaces.

Definition 1. Tensor space

We define a *tensor space* \mathbb{T} of rank r as a vector space such that its canonical basis is a cartesian product of the canonical bases of r other vector spaces. Its shape is denoted $n_1 \times n_2 \times \cdots \times n_r$, where the $\{n_k\}$ are the dimensions of the vector spaces.

Definition 2. Tensor product of vector spaces

Given r vector spaces $\mathbb{V}_1, \mathbb{V}_2, \dots, \mathbb{V}_r$, their *tensor product* is the tensor space \mathbb{T} spanned by the cartesian product of their canonical bases under coordinate-wise sum and outer product.

We use the notation $\mathbb{T} = \bigotimes_{k=1}^r \mathbb{V}_k$.

Remark. This simpler definition is indeed equivalent with the definition of the tensor product given in (Hackbusch, 2012, p. 51). The drawback of our definition is that it depends on the canonical bases, which at first can seem limiting as being canon implies that they are bounded to a certain system of coordinates. However this is not a concern in our domain as we need not distinguish tensors from their representation.

Remark. For naming convenience, from now on, we will distinguish between the terms *linear space* and *vector space* *i.e.* we will abusively use the term *vector space* only to refer to a linear space that is seen as a tensor space of rank 1. If we don't know its rank, we rather use the term *linear space*. We also make a clear distinction between the terms *dimension* (that is, for a tensor space it is equal to $\prod_{k=1}^r n_k$) and the term *rank* (equal to r). Note that some authors use the term *order* instead of *rank* (*e.g.* Hackbusch, 2012) as the latter is affected to another notion.

Definition 3. Tensor

A *tensor* t is an object of a tensor space. The *shape* of t , which is the same as the shape of the tensor space it belongs to, is denoted $n_1^{(t)} \times n_2^{(t)} \times \dots \times n_r^{(t)}$.

1.1.2 Manipulation

In this subsection, we describe notations and operators used to manipulate data stored in tensors.

Definition 4. Indexing

An *entry* of a tensor $t \in \mathbb{T}$ is one of its scalar coordinates in the canonical basis, denoted $t[i_1, i_2, \dots, i_r]$.

More precisely, if $\mathbb{T} = \bigotimes_{k=1}^r \mathbb{V}_k$, with bases $((e_k^i)_{i=1, \dots, n_k})_{k=1, \dots, r}$, then we have

$$t = \sum_{i_1=1}^{n_1} \cdots \sum_{i_r=1}^{n_r} t[i_1, i_2, \dots, i_r] (e_1^{i_1}, \dots, e_r^{i_r})$$

The cartesian product $\mathbb{I} = \prod_{k=1}^r \llbracket 1, n_k \rrbracket$ is called the *index space* of \mathbb{T}

Remark. When using an index i_k for an entry of a tensor t , we implicitly assume that $i_k \in \llbracket 1, n_k^{(t)} \rrbracket$ if nothing is specified.

Definition 5. Subtensor

A *subtensor* t' is a tensor of same rank composed of entries of t that are contiguous in the indexing, with at least one entry per rank. We denote $t' = t[l_1:u_1, l_2:u_2, \dots, l_r:u_r]$, where the $\{l_k\}$ and the $\{u_k\}$ are the lower and upper bounds of the indices used by the entries that compose t' .

Remark. We don't necessarily write the lower bound index if it is equal to 1, neither the upper bound index if it is equal to $n_k^{(t)}$.

Definition 6. Slicing

A *slice* operation, along the last ranks $\{r_1, r_2, \dots, r_s\}$, and indexed by $(i_{r_1}, i_{r_2}, \dots, i_{r_s})$,

is a morphism $s : \mathbb{T} = \bigotimes_{k=1}^r \mathbb{V}_k \rightarrow \bigotimes_{k=1}^{r-s} \mathbb{V}_k$, such that:

$$s(t)[i'_1, i'_2, \dots, i'_{r-s}] = t[i'_1, i'_2, \dots, i'_{r-s}, i_{r_1}, i_{r_2}, \dots, i_{r_s}]$$

$$i.e. \quad s(t) := t[:, :, \dots, :, i_{r_1}, i_{r_2}, \dots, i_{r_s}]$$

where $:=$ means that entries of the right operand are assigned to the left operand. We denote $t_{i_{r_1}, i_{r_2}, \dots, i_{r_s}}$ and call it the *slice* of t . Slicing along a

subset of ranks that are not the lasts is defined similarly. $s(\mathbb{T})$ is called a *slice subspace*.

Definition 7. Flattening

A *flatten* operation is an isomorphism $f : \mathbb{T} \rightarrow \mathbb{V}$, between a tensor space \mathbb{T} of rank r and an n -dimensional vector space \mathbb{V} , where $n = \prod_{k=1}^r n_k$. It is characterized by a bijection in the index spaces $g : \prod_{k=1}^r \llbracket 1, n_k \rrbracket \rightarrow \llbracket 1, n \rrbracket$ such that

$$\forall t \in \mathbb{T}, f(t)[g(i_1, i_2, \dots, i_r)] = f(t[i_1, i_2, \dots, i_r])$$

We call an inverse operation a *de-flatten* operation.

Remark. **Row major ordering**

The choice of g determines in which order the indexing is made. g is reminiscent of how data of multidimensional arrays or tensors are stored internally by programming languages. In most tensor manipulation languages, incrementing the memory address (*i.e.* the output of g) will first increment the last index i_r if $i_r < n_r$ (and if else $i_r = n_r$, then $i_r := 1$ and ranks are ordered in reverse lexicographic order to decide what indices are incremented). This is called *row major ordering*, as opposed to *column major ordering*. That is, in row major, g is defined as

$$g(i_1, i_2, \dots, i_r) = \sum_{p=1}^r \left(\prod_{k=p+1}^r n_k \right) i_p \quad (1.1)$$

Definition 8. Reshaping

A *reshape* operation is an isomorphism defined on a tensor space $\mathbb{T} = \bigotimes_{k=1}^r \mathbb{V}_k$ such that some of its basis vector spaces $\{\mathbb{V}_k\}$ are de-flattened and some of its slice subspaces are flattened.

1.1.3 Binary operations

We define binary operations on tensors that we'll later have use for. In particular, we define *tensor contraction* which is sometimes called *tensor multiplication*, *tensor product* or *tensor dotproduct* by other sources. We also define *convolution* and *pooling* which serve as the common building blocks of convolution neural network architectures (see Section 1.2.7).

Definition 9. Contraction

A *tensor contraction* between two tensors, along ranks of same dimensions, is defined by natural extension of the dot product operation to tensors.

More precisely, let \mathbb{T}_1 a tensor space of shape $n_1^{(1)} \times n_2^{(1)} \times \cdots \times n_{r_1}^{(1)}$, and \mathbb{T}_2 a tensor space of shape $n_1^{(2)} \times n_2^{(2)} \times \cdots \times n_{r_2}^{(2)}$, such that $\forall k \in \llbracket 1, s \rrbracket, n_{r_1-(s-k)}^{(1)} = n_k^{(2)}$, then the tensor contraction between $t_1 \in \mathbb{T}_1$ and $t_2 \in \mathbb{T}_2$ is defined as:

$$\left\{ \begin{array}{l} t_1 \otimes t_2 = t_3 \in \mathbb{T}_3 \text{ of shape } n_1^{(1)} \times \cdots \times n_{r_1-s}^{(1)} \times n_{s+1}^{(2)} \times \cdots \times n_{r_2}^{(2)} \text{ where} \\ t_3[i_1^{(1)}, \dots, i_{r_1-s}^{(1)}, i_{s+1}^{(2)}, \dots, i_{r_2}^{(2)}] = \\ \sum_{k_1=1}^{n_1^{(2)}} \cdots \sum_{k_s=1}^{n_s^{(2)}} t_1[i_1^{(1)}, \dots, i_{r_1-s}^{(1)}, k_1, \dots, k_s] t_2[k_1, \dots, k_s, i_{s+1}^{(2)}, \dots, i_{r_2}^{(2)}] \end{array} \right.$$

For the sake of simplicity, we omit the case where the contracted ranks are not the last ones for t_1 and the first ones for t_2 . But this definition still holds in the general case subject to a permutation of the indices.

Definition 10. Covariant and contravariant indices

Given a tensor contraction $t_1 \otimes t_2$, indices of the left hand operand t_1 that are not contracted are called *covariant* indices. Those that are contracted are called *contravariant* indices. For the right operand t_2 , the naming convention is the opposite. The set of covariant and contravariant indices of both operands are called the *transformation laws* of the tensor contraction.

Remark. Transformation law independency

Contrary to most mathematical definitions, tensors in deep learning are inde-

pendent of any transformation law, so that they must be specified for tensor contractions.

Remark. **Einstein summation convention**

Using subscript notation for covariant indices and superscript notation for contravariant indices, the previous tensor contraction can be written using the Einstein summation convention as:

$$t_{1,i_1^{(1)} \dots i_{r_1-s}^{(1)}}{}^{k_1 \dots k_s} t_{2,k_1 \dots k_s}{}^{i_{s+1}^{(2)} \dots i_{r_2}^{(2)}} = t_{3,i_1^{(1)} \dots i_{r_1-s}^{(1)}}{}^{i_{s+1}^{(2)} \dots i_{r_2}^{(2)}} \quad (1.2)$$

Dot product $u_k v^k = \lambda$ and matrix product $A_i{}^k B_k{}^j = C_i{}^j$ are common examples of tensor contractions.

Proposition 11. A contraction can be rewritten as a matrix product.

Proof. Using notation of (1.2), with the reshapings $t_1 \mapsto T_1$, $t_2 \mapsto T_2$ and $t_3 \mapsto T_3$ defined by grouping all covariant indices into a single index and all contravariant indices into another single index, we can rewrite

$$T_{1,g_i(i_1^{(1)}, \dots, i_{r_1-s}^{(1)})}{}^{g_k(k_1, \dots, k_s)} T_{2,g_k(k_1, \dots, k_s)}{}^{g_j(i_{s+1}^{(2)}, \dots, i_{r_2}^{(2)})} = T_{3,g_i(i_1^{(1)}, \dots, i_{r_1-s}^{(1)})}{}^{g_j(i_{s+1}^{(2)}, \dots, i_{r_2}^{(2)})}$$

where g_i , g_k and g_j are bijections defined similarly as in (1.1). \square

Definition 12. Convolution

The n -dimensional convolution, denoted $*^n$, between $t_1 \in \mathbb{T}_1$ and $t_2 \in \mathbb{T}_2$, where \mathbb{T}_1 and \mathbb{T}_2 are of the same rank n such that $\forall p \in \llbracket 1, n \rrbracket, n_p^{(1)} \geq n_p^{(2)}$, is defined as:

$$\left\{ \begin{array}{l} t_1 *^n t_2 = t_3 \in \mathbb{T}_3 \text{ of shape } n_1^{(3)} \times \dots \times n_n^{(3)} \text{ where} \\ \forall p \in \llbracket 1, n \rrbracket, n_p^{(3)} = n_p^{(1)} - n_p^{(2)} + 1 \\ t_3[i_1, \dots, i_n] = \sum_{k_1=1}^{n_1^{(2)}} \dots \sum_{k_n=1}^{n_n^{(2)}} t_1[i_1 + n_1^{(2)} - k_1, \dots, i_n + n_n^{(2)} - k_n] t_2[k_1, \dots, k_n] \end{array} \right.$$

Proposition 13. A convolution can be rewritten as a matrix product.

Proof. Let $t_1 *^n t_2 = t_3$ defined as previously with $\mathbb{T}_1 = \bigotimes_{k=1}^r \mathbb{V}_k^{(1)}$, $\mathbb{T}_2 = \bigotimes_{k=1}^r \mathbb{V}_k^{(2)}$. Let $t'_1 \in \bigotimes_{k=1}^r \mathbb{V}_k^{(1)} \otimes \bigotimes_{k=1}^r \mathbb{V}_k^{(2)}$ such that $t'_1[i_1, \dots, i_n, k_1, \dots, k_n] = t_1[i_1 + n_1^{(2)} - k_1, \dots, i_n + n_n^{(2)} - k_n]$, then

$$t_3[i_1, \dots, i_n] = \sum_{k_1=1}^{n_1^{(2)}} \cdots \sum_{k_n=1}^{n_n^{(2)}} t'_1[i_1, \dots, i_n, k_1, \dots, k_n] t_2[k_1, \dots, k_n]$$

where we recognize a tensor contraction. Proposition 11 concludes. \square

The two following operations are meant to further decrease the shape of the resulting output.

Definition 14. Strided convolution

The n -dimensional *strided* convolution, with strides $s = (s_1, s_2, \dots, s_n)$, denoted $*_s^n$, between $t_1 \in \mathbb{T}_1$ and $t_2 \in \mathbb{T}_2$, where \mathbb{T}_1 and \mathbb{T}_2 are of the same rank n such that $\forall p \in \llbracket 1, n \rrbracket, n_p^{(1)} \geq n_p^{(2)}$, is defined as:

$$\begin{cases} t_1 *_s^n t_2 = t_4 \in \mathbb{T}_4 \text{ of shape } n_1^{(4)} \times \cdots \times n_n^{(4)} \text{ where} \\ \forall p \in \llbracket 1, n \rrbracket, n_p^{(4)} = \lfloor \frac{n_p^{(1)} - n_p^{(2)} + 1}{s_p} \rfloor \\ t_4[i_1, \dots, i_n] = (t_1 *_n t_2)[(i_1 - 1)s_1 + 1, \dots, (i_n - 1)s_n + 1] \end{cases}$$

Remark. Informally, a strided convolution is defined as if it were a regular subsampling of a convolution. They match if $s = (1, 1, \dots, 1)$.

Definition 15. Pooling

Let a real-valued function f defined on all tensor spaces of any shape, *e.g.* the *max* or *average* function. An f -pooling operation is a mapping $t \mapsto t'$ such that each entry of t' is an image by f of a subtensor of t .

Remark. Usually, the set of subtensors that are reduced by f into entries of t' are defined by a regular partition of the entries of t .

1.2 Neural Networks

A feed-forward neural network could originally be formalized as a composite function chaining linear and non-linear functions (Rumelhart et al., 1985; LeCun et al., 1989; LeCun and Bengio, 1995), even up until the important breakthroughs that generated a surge of interest in the field (Hinton et al., 2012; Krizhevsky et al., 2012; Simonyan and Zisserman, 2014). However, in more recent advances, more complex architectures have emerged (Szegedy et al., 2015; He et al., 2016; Zoph and Le, 2016; Huang et al., 2017), such that the former formalization does not suffice. We provide a definition for the first kind of neural networks (Definition 16) and use it to present its related concepts. Then we give a more generic definition (Definition 20).

Note that in this manuscript, we only consider neural networks that are *feed-forward* (Zell, 1994; Wikipedia, 2018b).

1.2.1 Simple formalization

We denote by I_f the *domain of definition* of a function f ("I" stands for "input") and by $O_f = f(I_f)$ its *image* ("O" stands for "output"), and we represent it as $I_f \xrightarrow{f} O_f$.

Definition 16. Neural network (simply connected)

Let f be a function such that I_f and O_f are vector or tensor spaces.

f is a (*simply connected*) *neural network function* if there are a series of affine functions $(g_k)_{k=1,2,\dots,L}$ and a series of non-linear derivable univariate functions $(h_k)_{k=1,2,\dots,L}$ such that:

$$\begin{cases} \forall k \in \llbracket 1, L \rrbracket, f_k = h_k \circ g_k, \\ I_f = I_{f_1} \xrightarrow{f_1} O_{f_1} \cong I_{f_2} \xrightarrow{f_2} \dots \xrightarrow{f_L} O_{f_L} = O_f, \\ f = f_L \circ \dots \circ f_2 \circ f_1 \end{cases}$$

The couple (g_k, h_k) is called the k -th *layer* of the neural network. L is its

depth. For $x \in I_f$, we denote by $x_k = f_k \circ \dots \circ f_2 \circ f_1(x)$ the *activations* of the k -th layer. We denote by \mathcal{N} the set of neural network functions.

Definition 17. Activation function

An *activation function* h is a real-valued univariate function that is non-linear and derivable, that is also defined by extension on any linear space with the functional notation $h(v)[i] = h(v[i])$.

Definition 18. Layer

A couple (g, h) , where g is an affine or linear function, and h is an activation function is called a *layer*. The set of layers is denoted \mathcal{L} .

Remark. **Adoption of ReLU activations**

Historically, sigmoidal and tanh activations were mostly used (Cybenko, 1989; LeCun et al., 1989). However in recent practice, the *rectified linear unit* (ReLU), which implements the *rectifier* function $h : x \mapsto \max(0, x)$ with convention $h'(0) = 0$ (first introduced as the *positive part*, Jarrett et al., 2009), is the most used activation, as it was demonstrated to be faster and to obtain better results (Glorot et al., 2011). ReLU originated numerous variants *e.g.* *leaky rectified linear unit* (Maas et al., 2013), *parametric rectified linear unit* (PReLU, He et al., 2015), *exponential linear unit* (ELU, Clevert et al., 2015), *scaled exponential linear unit* (SELU, Klambauer et al., 2017).

Remark. **Universal approximation**

Early researches have shown that neural networks with one level of depth can approximate any real-valued function defined on a compact subset of \mathbb{R}^n . This result was first proved for sigmoidal activations (Cybenko, 1989), and then it was shown it did not depend on the sigmoidal activations (Hornik et al., 1989; Hornik, 1991).

For example, for the application of supervised learning when a neural network is trained from data (see Section 1.2.4), this result is quite important because it brings theoretical justification that the objective exists (even though it doesn't inform whether an algorithm to approach it exists or is efficient).

Remark. Computational difficulty

However, reaching such objective is a computationally difficult problem, which drove back interest from the field. Thanks to better hardware and to using better initialization schemes that speed up learning, researchers started to report more successes with deep neural networks (Hinton et al., 2006; Glorot and Bengio, 2010) ; see (Bengio, 2009) for a review of this period. It ultimately came to a surge of interest in the field after a significant breakthrough on the imagenet dataset (Deng et al., 2009) with a deep convolutional architecture (Krizhevsky et al., 2012), see Section 1.2.7. The use of the fast ReLU activation function (Glorot et al., 2011) as well as leveraging graphical processing units with CUDA (Nickolls et al., 2008) were also key factors in overcoming this computational difficulty.

Remark. Expressivity and expressive efficiency

The study of the *expressivity* (also called *representational power*) of families of neural networks is the field that is interested in the range of functions that can be realized or approximated by this family (Håstad and Goldmann, 1991; Pascanu et al., 2013). In general, given a maximal error ϵ and an objective F , the more expressive is a family $N \subset \mathcal{N}$, the more likely it is to contain an approximation $f \in N$ such that $d(f, F) < \epsilon$. However, if we consider the approximation $f_{min} \in N$ that have the lowest number of neurons, it is possible that f_{min} is still too large and may be unpractical. For this reason, expressivity is often studied along the related notion of *expressive efficiency* (Delalleau and Bengio, 2011; Cohen et al., 2018a).

Remark. Rectifier neural networks

Of particular interest for the intuition is a result stating that a simply connected neural networks with only ReLU activations (a rectifier neural network) is a piecewise linear function (Pascanu et al., 2013; Montufar et al., 2014), and that conversely any piecewise linear function is also a rectifier neural network such that an upper bound of its depth is logarithmically related to the input dimension (Arora et al., 2018, th. 2.1.). Their expressive

efficiency have also been demonstrated compared to neural networks using threshold or sigmoid activations (Pan and Srikumar, 2016).

Remark. Benefits of depth

Expressive efficiency analysis have demonstrated the benefits of depth, *i.e.* a shallow neural network would need an unfeasible large number of neurons to approximate the function of a deep neural network (*e.g.* Delalleau and Bengio, 2011; Bianchini and Scarselli, 2014; Poggio et al., 2015; Eldan and Shamir, 2016; Poole et al., 2016; Raghu et al., 2016; Cohen and Shashua, 2016; Mhaskar et al., 2016; Lin et al., 2017; Arora et al., 2018).

Remark. Bias

Note that affine functions \tilde{g} can be written as a sum between a linear function g and a constant vector b which is called the *bias*. It augments the expressivity of the neural network's family of functions. For notational convenience, we will often omit to write down the biases in the layer's equations.

1.2.2 Generic formalization

The former neural networks are said to be *simply connected* because each layer only takes as input the output of the previous one. We'll give a more general definition after first defining branching operations.

Definition 19. Branching

A *binary branching operation* between two tensors, $x_{k_1} \bowtie x_{k_2}$, outputs, subject to shape compatibility, either their addition, either their concatenation along a rank, or their concatenation as a list.

A *branching operation* between n tensors, $x_{k_1} \bowtie x_{k_2} \bowtie \cdots \bowtie x_{k_n}$, is a composition of binary branching operations, or is the identity function Id if $n = 1$. Branching operations are also naturally defined on tensor-valued functions through their realizations.

Definition 20. Neural network (generic definition)

The set of *neural network* functions \mathcal{N} is defined inductively as follows

1. $Id \in \mathcal{N}$
2. $f \in \mathcal{N} \wedge (g, h) \in \mathcal{L} \wedge O_f \subset I_g \Rightarrow h \circ g \circ f \in \mathcal{N}$
3. for all shape compatible branching operations:
 $f_1, f_2, \dots, f_n \in \mathcal{N} \Rightarrow f_1 \bowtie f_2 \bowtie \dots \bowtie f_n \in \mathcal{N}$

Remark. **Examples**

The neural network proposed in (Szegedy et al., 2015), called *Inception*, use depth-wise concatenation of feature maps. Residual networks (ResNets, He et al., 2016) make use of *residual connections*, also called *skip connections*, *i.e.* an activation that is used as input in a lower level is added to another activation at an upper level. Densely connected networks (DenseNets, Huang et al., 2017) have their activations concatenated with all lower level activations. These neural networks had demonstrated state of the art performances on the imagenet classification challenge (Deng et al., 2009), outperforming simply connected neural networks.

Remark. **Benefits of branching operations**

Recent works have provided rationales supporting benefits of using branching operations, thus giving justifications for architectures obtained with the generic formalization. In particular, (Cohen et al., 2018a) have analyzed the impact of residual connections used in Wavenet-like architectures (Van Den Oord et al., 2016) in terms of expressive efficiency (see Remark 1.2.1) using tools from the field of tensor analysis ; (Orhan and Pitkow, 2018) have empirically demonstrated that skip connections can resolve some inefficiency problems inherent of fully-connected networks (dead activations, activations that are always equal, linearly dependent sets of activations).

For layer indexing convenience, we still use the simple formalization in the subsequent subsections, even though the presentation would be similar with the generic formalization.

1.2.3 Interpretation

Until now, we have formally introduced a neural network as a mathematical function. As its name suggests, such function can be interpreted from a connectivity viewpoint (LeCun, 1987).

Definition 21. Connectivity matrix

Let g a linear function. Without loss of generality subject to a flattening, let's suppose I_g and O_g are vector spaces. Then there exists a *connectivity matrix* W_g , such that:

$$\forall x \in I_g, g(x) = W_g x$$

We denote W_k the connectivity matrix of the k -th layer.

Remark. **Biological inspiration**

A (*computational*) *neuron* is a computational unit that is biologically inspired (McCulloch and Pitts, 1943). Each neuron is capable of:

1. receiving modulated signals from other neurons and aggregate them,
2. applying to the result a derivable activation,
3. passing the signal to other neurons.

That is to say, each domain $\{I_{f_k}\}$ and O_f can be interpreted as a layer of neurons, with one neuron for each dimension. The connectivity matrices $\{W_k\}$ describe the connections between each successive layers. A neuron is illustrated on Figure 1.1.

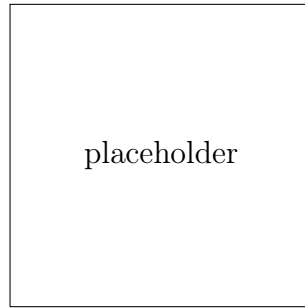


Figure 1.1: A neuron

1.2.4 Parameterization and training

TODO: pass on this section

Given an objective function F , training is the process of incrementally modifying a neural network f upon obtaining a better approximation of F . The most used training algorithms are based on gradient descent, as proposed in (Widrow and Hoff, 1960). These algorithms became popular since (Rumelhart et al., 1985). Informally, f is parameterized with initial weights that characterize its linear parts. These weights are modified step by step in the opposite direction of their gradient with respect to a loss. All possible realizations of f through its weights draw a family \mathcal{N} whose expressivity is crucial for the success of the training. The common points between f and other objects of \mathcal{N} define what is called a neural network *architecture*. That is

We present gradient based learning more formally in what follows.

Definition 22. Architecture Let $f \in \mathcal{N}$ with weights $(\theta_k)_k \in$.

Remark. Gradient descent

The most used training algorithms are based on gradient descent, as proposed in (Widrow and Hoff, 1960). These algorithms became popular since (Rumelhart et al., 1985). In order to be trained, f is parameterized with

initial weights that characterize its linear parts. These weights are modified step by step in the opposite direction of their gradient with respect to a loss.

Remark. Architecture

All possible realizations of f through its weights draw a family N which expressivity is crucial for the success of the training. The common points between f and other objects of N define what is called a neural network *architecture*.

Definition 23. Weights

Let consider the k -th layer of a neural network f . We define its weights as coordinates of a vector θ_k , called the *weight kernel*, such that:

$$\forall(i, j), \begin{cases} \exists p, W_k[i, j] := \theta_k[p] \\ \text{or } W_k[i, j] = 0 \end{cases}$$

A weight p that appears multiple times in W_k is said to be *shared*. Two parameters of W_k that share a same weight p are said to be *tied*. The number of weights of the k -th layer is $n_1^{(\theta_k)}$.

Remark. Learning

A *loss* function \mathcal{L} penalizes the output $x_L = f(x)$ relatively to the approximation error $|f(x) - F(x)|$. Gradient w.r.t. θ_k , denoted $\vec{\nabla}_{\theta_k}$, is used to update the weights via an optimization algorithm based on gradient descent and a learning rate α , that is:

$$\theta_k^{(\text{new})} = \theta_k^{(\text{old})} - \alpha \cdot \vec{\nabla}_{\theta_k} \left(\mathcal{L} \left(x_L, \theta_k^{(\text{old})} \right) + R \left(\theta_k^{(\text{old})} \right) \right) \quad (1.3)$$

where α can be a scalar or a vector, \cdot can denote outer or pointwise product, and R is a regularizer. They depend on the optimization algorithm.

TODO: examples of optimization

Remark. **Linear complexity**

The complexity of computing the gradients is linear with the number of weights.

Proof. Without loss of generality, we assume that the neural network is simply connected. Thanks to the chain rule, $\vec{\nabla}_{\theta_k}$ can be computed using gradients that are w.r.t. x_k , denoted $\vec{\nabla}_{x_k}$, which in turn can be computed using gradients w.r.t. outputs of the next layer $k+1$, up to the gradients given on the output layer.

That is:

$$\vec{\nabla}_{\theta_k} = J_{\theta_k}(x_k) \vec{\nabla}_{x_k} \quad (1.4)$$

$$\begin{aligned} \vec{\nabla}_{x_k} &= J_{x_k}(x_{k+1}) \vec{\nabla}_{x_{k+1}} \\ \vec{\nabla}_{x_{k+1}} &= J_{x_{k+1}}(x_{k+2}) \vec{\nabla}_{x_{k+2}} \end{aligned} \quad (1.5)$$

$$\begin{aligned} &\dots \\ \vec{\nabla}_{x_{L-1}} &= J_{x_{L-1}}(x_L) \vec{\nabla}_{x_L} \end{aligned}$$

Obtaining,

$$\vec{\nabla}_{\theta_k} = J_{\theta_k}(x_k) \left(\prod_{p=k}^{L-1} J_{x_p}(x_{p+1}) \right) \vec{\nabla}_{x_L} \quad (1.6)$$

where $J_{\text{wrt}}(.)$ are the respective jacobians which can be determined with the layer's expressions and the $\{x_k\}$; and $\vec{\nabla}_{x_L}$ can be determined using \mathcal{L} , R and x_L . \square

This allows to compute the gradients with a complexity that is linear with the number of weights (only one computation of the activations), instead of being quadratic if it were done with the difference quotient expression of the derivatives (one more computation of the activations for each weight).

Remark. Backpropagation

We can remark that (1.5) rewrites as

$$\begin{aligned}\vec{\nabla}_{x_k} &= J_{x_k}(x_{k+1}) \vec{\nabla}_{x_{k+1}} \\ &= J_{x'_k}(h(x'_k)) J_{x_k}(W_k x_k) \vec{\nabla}_{x_{k+1}}\end{aligned}\tag{1.7}$$

where $x'_k = W_k x_k$, and these jacobians can be expressed as:

$$J_{x'_k}(h(x'_k))[i, j] = \delta_i^j h'(x'_k[i])\tag{1.8}$$

$$J_{x'_k}(h(x'_k)) = I h'(x'_k)$$

$$J_{x_k}(W_k x_k) = W_k^T\tag{1.9}$$

That means that we can write $\vec{\nabla}_{x_k} = (\tilde{h}_k \circ \tilde{g}_k)(\vec{\nabla}_{x_{k+1}})$ such that the connectivity matrix \widetilde{W}_k is obtained by transposition. This can be interpreted as gradient calculation being a *back-propagation* on the same neural network, in opposition of the *forward-propagation* done to compute the output.

TODO: Overfitting remark

1.2.5 Examples of layer

Definition 24. Connections

The set of *connections* of a layer (g, h) , denoted C_g , is defined as:

$$C_g = \{(i, j), \exists p, W_g[i, j] := \theta_g[p]\}$$

We have $0 \leq |C_g| \leq n_1^{(W_g)} n_2^{(W_g)}$.

Definition 25. Dense layer

A *dense layer* (g, h) is a layer such that $|C_g| = n_1^{(W_g)} n_2^{(W_g)}$, *i.e.* all possible connections exist. The map $(i, j) \mapsto p$ is usually a bijection, meaning that

there is no weight sharing.

Definition 26. Partially connected layer

A *partially connected layer* (g, h) is a layer such that $|C_g| < n_1^{(W_g)} n_2^{(W_g)}$.

A *sparse connected layer* (g, h) is a layer such that $|C_g| \ll n_1^{(W_g)} n_2^{(W_g)}$.

Definition 27. Convolutional layer

A n -dimensional convolutional layer (g, h) is such that the weight kernel θ_g can be reshaped into a tensor w of rank $n + 2$, and such that

$$\begin{cases} I_g \text{ and } O_g \text{ are tensor spaces of rank } n + 1 \\ \forall x \in I_g, g(x) = (g(x))_q = \sum_p x_p *^n w_{p,q} \forall q \end{cases}$$

where p and q index slices along the last ranks.

Definition 28. Feature maps and input channels

The slices $g(x)_q$ are typically called *feature maps*, and the slices x_p are called *input channels*. Let's denote by $n_o = n_{n+1}^{(O_g)}$ and $n_i = n_{n+1}^{(I_g)}$ the number of feature maps and input channels. In other words, Definition 27 means that for each feature maps, a convolution layer computes n_i convolutions and sums them, computing a total of $n_i \times n_o$ convolutions.

Remark. Note that because they are simply summed, entries of two different input channels that have the same coordinates are assumed to share some sort of relationship. For instance on images, entries of each input channel (typically corresponding to Red, Green and Blue channels) that have the same coordinates share the same pixel location.

Remark. Given a tensor input x , the n -dimensional convolutions between the inputs channels x_p and slices of a weight tensor $w_{p,q}$ would result in outputs y_q of shape $n_1^{(x)} - n_1^{(w)} + 1 \times \dots \times n_n^{(x)} - n_n^{(w)} + 1$. So, in order to preserve shapes, a padding operation must pad x with $n_1^{(w)} - 1 \times \dots \times n_n^{(w)} - 1$ zeros beforehand. For example, the padding function of the library *tensorflow* (Abadi et al.,

2015) pads each rank with a balanced number of zeros on the left and right indices (except if $n_t^{(w)} - 1$ is odd then there is one more zero on the left).

Definition 29. Padding

A convolutional layer with *padding* (g, h) is such that g can be decomposed as $g = g_{\text{pad}} \circ g'$, where g' is the linear part of a convolution layer as in Definition 27, and g_{pad} is an operation that pads zeros to its inputs such that g preserves tensor shapes.

Remark. One asset of padding operations is that they limit the possible loss of information on the borders of the subsequent convolutions, as well as preventing a decrease in size. Moreover, preserving shape is needed to build some neural network architectures, especially for ones with branching operations *e.g.* Remark 1.2.2. On the other hand, they increase memory and computational footprints.

Proposition 30. Connectivity matrix of a convolution with padding

A convolutional layer with padding (g, h) is equivalently defined as W_g being a $n_i \times n_o$ block matrix such that its blocks are Toeplitz matrices.

Proof. Let's consider the slices indexed by p and q , and to simplify the notations, let's drop the subscripts p, q . We recall from Definition 12 that

$$\begin{aligned}
 y &= (x *^n w)[j_1, \dots, j_n] \\
 &= \sum_{k_1=1}^{n_1^{(w)}} \cdots \sum_{k_n=1}^{n_n^{(w)}} x[j_1 + n_1^{(w)} - k_1, \dots, j_n + n_n^{(w)} - k_n] w[k_1, \dots, k_n] \\
 &= \sum_{i_1=j_1}^{j_1+n_1^{(w)}-1} \cdots \sum_{i_n=j_n}^{j_n+n_n^{(w)}-1} x[i_1, \dots, i_n] w[j_1 + n_1^{(w)} - i_1, \dots, j_n + n_n^{(w)} - i_n] \\
 &= \sum_{i_1=1}^{n_1^{(x)}} \cdots \sum_{i_n=1}^{n_n^{(x)}} x[i_1, \dots, i_n] \tilde{w}[i_1, j_1, \dots, i_n, j_n] \\
 &\text{where } \tilde{w}[i_1, j_1, \dots, i_n, j_n] =
 \end{aligned}$$

$$\begin{cases} w[j_1 + n_1^{(w)} - i_1, \dots, j_n + n_n^{(w)} - i_n] & \text{if } \forall t, 0 \leq i_t - j_t \leq n_t^{(w)} - 1 \\ 0 & \text{otherwise} \end{cases}$$

Using Einstein summation convention as in (1.2) and permuting indices, we recognize the following tensor contraction

$$y_{j_1 \dots j_n} = x_{i_1 \dots i_n} \tilde{w}^{i_1 \dots i_n}_{j_1 \dots j_n} \quad (1.10)$$

Following Proposition 11, we reshape (1.10) as a matrix product. To reshape $y \mapsto Y$, we use the row major order bijections g_j as in (1.1) defined onto $\{(j_1, \dots, j_n), \forall t, 1 \leq j_t \leq n_t^{(y)}\}$. To reshape $x \mapsto X$, we use the same row major order bijection g_j , however defined on the indices that support non zero-padded values, so that zero-padded values are lost after reshaping. That is, we use a bijection g_i such that $g_i(i_1, i_2, \dots, i_n) = g_j(i_1 - o_1, i_2 - o_2, \dots, i_n - o_n)$ defined if and only if $\forall t, 1 + o_t \leq i_t \leq n_t^{(y)}$, where the $\{o_t\}$ are the starting offsets of the non zero-padded values. $\tilde{w} \mapsto W$ is reshaped by using g_j for its covariant indices, and g_i for its contravariant indices. The entries lost by using g_i do not matter because they would have been nullified by the resulting matrix product. We remark that W is exactly the block (p, q) of W_g (and not of $W_{g'}$). Now let's prove that it is a Toeplitz matrix.

Thanks to the linearity of the expression (1.1) of g_j , by denoting $i'_t = i_t - o_t$, we obtain

$$g_i(i_1, i_2, \dots, i_n) - g_j(j_1, j_2, \dots, j_n) = g_j(i'_1 - j_1, i'_2 - j_2, \dots, i'_n - j_n) \quad (1.11)$$

To simplify the notations, let's drop the arguments of g_i and g_j . By bijectivity of g_j , (1.11) tells us that $g_i - g_j$ remains constant if and only if $i'_t - j_t$ remains

constant for all t . Recall that

$$W[g_i, g_j] = \begin{cases} w[j_1 + n_1^{(w)} - i'_1, \dots, j_n + n_n^{(w)} - i'_n] & \text{if } \forall t, 0 \leq i'_t - j_t \leq n_t^{(w)} - 1 \\ 0 & \text{otherwise} \end{cases} \quad (1.12)$$

Hence, on a diagonal of W , $g_i - g_j$ remaining constant means that $W[g_i, g_j]$ also remains constants. So W is a Toeplitz matrix.

The converse is also true as we used invertible functions in the index spaces through the proof. \square

Remark. The former proof makes clear that the result doesn't hold in case there is no padding. This is due to border effects when the index of the n^{th} rank resets in the definition of the row-major ordering function g_j that would be used. Indeed, under appropriate definitions, the matrices could be seen as almost Toeplitz.

Remark. Comparatively with dense layers, convolution layers enjoy a significant decrease in the number of weights. For example, an input 2×2 convolution on images with 3-color input channels, would breed only 12 weights per feature maps, independently of the numbers of input neurons. On image datasets, their usage also breeds a significant boost in performance compared with dense layers (Krizhevsky et al., 2012), for they allow to take advantage of the topology of the inputs while dense layers don't (LeCun and Bengio, 1995). A more thorough comparison and explanation of their assets will be discussed in Section ??.

Definition 31. Stride

A convolutional layer with *stride* is a convolutional layer that computes strided convolutions (with $\text{stride} > 1$) instead of convolutions.

Definition 32. Pooling

A layer with *pooling* (g, h) is such that g can be decomposed as $g = g' \circ g_{\text{pool}}$, where g_{pool} is a pooling operation.

Remark. Downscaling

Layers with stride or pooling downscale the signals that passes through the layer. These types of layers allows to compute features at a coarser level, giving the intuition that the deeper a layer is in the network, the more abstract is the information captured by the weights of the layer.

TODO: below

1.2.6 Examples of regularization

Remark. Overfitting TODO:

A layer with *dropout* (g, h) is such that $h = h_1 \circ h_2$, where (g, h_2) is a layer and h_1 is a dropout operation (Srivastava et al., 2014). When dropout is used, a certain number of neurons are randomly set to zero during the training phase, compensated at test time by scaling down the whole layer. This is done to prevent overfitting.

1.2.7 Examples of architecture

TODO: rephrase

A multilayer perceptron (MLP) (Hornik et al., 1989) is a neural network composed of only dense layers. A convolutional neural network (CNN) (LeCun et al., 1998) is a neural network composed of convolutional layers.

Neural networks are commonly used for machine learning tasks. For example, to perform supervised classification, we usually add a dense output layer $s = (g_{L+1}, h_{L+1})$ with as many neurons as classes. We measure the error between an output and its expected output with a discriminative loss function \mathcal{L} . During the training phase, the weights of the network are adapted for the classification task based on the errors that are back-propagated (Hornik et

al., 1989) via the chain rule and according to a chosen optimization algorithm (*e.g.* Bottou, 2010).

1.3 Graphs and signals

1.3.1 Basic definitions

1.3.1.1 Graphs

Definition 33. Graph

A *graph* G is defined as a couple of sets $\langle V, E \rangle$ where V is the set of *vertices*, also called *nodes*, and $E \subseteq \binom{V}{2}$ is the set of *edges*. For all $u, v \in V$ we define the relation $u \sim v \Leftrightarrow \{u, v\} \in E$. Unless stated otherwise, we will consider only *weighted* graphs *i.e.* each graph G is associated with a weight mapping $w : E \rightarrow \mathbb{R}^*$.

Figure 1.2 illustrates an example of a graph. Note that we employ interchangeably the terms *vertex* and *node*.

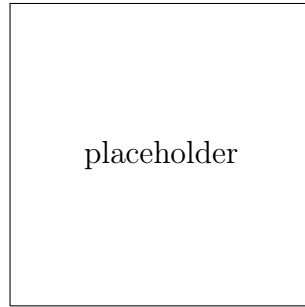


Figure 1.2: Example of a graph

Remark. According to this definition, we consider simple graphs *i.e.* no two edges share the same set of vertices and there is no self-loop.

Definition 34. Path

A *path* of length $n \in \mathbb{N}$ in a graph $G = \langle V, E \rangle$ is a sequence $(v_1 v_2 \cdots v_n)$ in V such that $\forall i, v_i \sim v_{i+1}$.

Remark. Our definition of graphs admit no self-loop so $\forall i, v_i \neq v_{i+1}$

Definition 35. Order

The order of a graph $G = \langle V, E \rangle$ is define as $\text{order}(G) = |V| \in \mathbb{N} \cup \{+\infty\}$

Definition 36. Adjacency matrix

The *adjacency matrix* of a finite graph $G = \langle V, E \rangle$ of order n , is a $n \times n$ real-valued matrix A associated to an indexing of $V = \{v_1, v_2, \dots, v_n\}$, such that

$$A[i, j] = \begin{cases} w(\{v_i, v_j\}) & \text{if } v_i \sim v_j \\ 0 & \text{otherwise} \end{cases}$$

Definition 37. Degree

The *degree* of a vertex $v \in V$ of a graph $G = \langle V, E \rangle$ is defined as $\deg(v) = |\{u \in V, u \sim v\}| \in \mathbb{N} \cup \{+\infty\}$.

The *degree* of the graph G is defined as $\deg(G) = \max_{v \in V} \{\deg(v)\}$.

A graph is said to be *regular* if \deg is constant on the vertices.

Definition 38. Degree matrix

The *degree matrix* of a finite graph $G = \langle V, E \rangle$ of order n , is the diagonal matrix D , associated to an indexing of $V = \{v_1, v_2, \dots, v_n\}$, such that $D = \text{diag}(\deg(v_1), \deg(v_2), \dots, \deg(v_n))$.

Definition 39. Laplacian matrix

The *laplacian matrix* of a graph $G = \langle V, E \rangle$ of order n , associated to an indexing of $V = \{v_1, v_2, \dots, v_n\}$, is defined as $L = D - A$, where D is the degree matrix and A is the adjacency matrix.

Definition 40. Digraph

A digraph is an oriented graph *i.e.* $E \subseteq V \times V - \{(v, v), v \in V\}$. Contrary to a graph, the weight mapping w , the relation \sim , the adjacency matrix A , and the laplacian matrix L are not symmetric. Notions defined on graphs naturally extends to digraphs where possible.

Definition 41. Bipartite graph

A *bipartite graph* is a triplet of sets $\langle V^{(1)}, V^{(2)}, E \rangle$, where $V^{(1)}$ and $V^{(2)}$ are sets of vertices, $V^{(1)} \cap V^{(2)} \neq \emptyset$, and $E \subseteq V^{(1)} \times V^{(2)}$. It is associated with a weight mapping $w : E \rightarrow \mathbb{R}^*$. Its adjacency matrix A is associated to indexings of $V^{(1)} = \{v_1^{(1)}, v_2^{(1)}, \dots, v_n^{(1)}\}$ and $V^{(2)} = \{v_1^{(2)}, v_2^{(2)}, \dots, v_n^{(2)}\}$, such that

$$A[i, j] = \begin{cases} w\left((v_i^{(1)}, v_j^{(2)})\right) & \text{if } (v_i^{(1)}, v_j^{(2)}) \in E \\ 0 & \text{otherwise} \end{cases}$$

Definition 42. Induced subgraph

The *subgraph* $\tilde{G} = \langle \tilde{V}, \tilde{E} \rangle$ of a graph $G = \langle V, E \rangle$, *induced* by $\tilde{V} \subseteq V$, is such that $\forall (u, v) \in \tilde{V}^2, u \stackrel{\tilde{G}}{\sim} v \Leftrightarrow u \stackrel{G}{\sim} v$.

Definition 43. Grid graph

A *grid graph* $G = \langle V, E \rangle$ is such that $V \cong \mathbb{Z}^2$, $v_1 \sim v_2 \Rightarrow \|v_2 - v_1\|_\infty \in \{0, 1\}$ and either one of the following is true:

$$\begin{cases} (i_1, j_1) \sim (i_2, j_2) \Leftrightarrow |i_2 - i_1| \text{ XOR } |j_2 - j_1| & (4 \text{ neighbours}) \\ (i_1, j_1) \sim (i_2, j_2) \Leftrightarrow |i_2 - i_1| \text{ AND } |j_2 - j_1| & (4 \text{ neighbours}) \\ (i_1, j_1) \sim (i_2, j_2) \Leftrightarrow |i_2 - i_1| \text{ OR } |j_2 - j_1| & (8 \text{ neighbours}) \end{cases}$$

A (*rectangular*) *grid graph* of size $n \times m$ is the subgraph of a grid graph induced by $\llbracket 1, n \rrbracket \times \llbracket 1, m \rrbracket$.

A *square grid graph* is a rectangular grid graph of square size.

1.3.1.2 (Real-valued) Signals**Definition 44. Signal space**

The *signal space* $\mathcal{S}(V)$, over the set V , is the linear space of real-valued functions defined on V .

We have $\dim(\mathcal{S}(V)) = |V| \in \mathbb{N} \cup \{+\infty\}$.

Remark. In particular, a vector space, and more generally a tensor space, are finite-dimensional signal spaces over any of their bases.

Definition 45. Signal

A *signal* over V , $s \in \mathcal{S}(V)$, is a function $s : V \rightarrow \mathbb{R}$.

An *entry* of a signal s is an image by s of some $v \in V$ and we denote $s[v]$.

If v is represented by a n -tuple, we can also write $s[v_1, v_2, \dots, v_n]$.

The *support* of a signal $s \in \mathcal{S}(V)$ is $\text{supp}(s) = \{v \in V, s[v] \neq 0\}$.

Definition 46. Graph signal

A *graph signal* over G is a signal over its vertex set. We denote by $\mathcal{S}(G)$ the graph signal space.

We have $\dim(\mathcal{S}(G)) = \text{order}(G) \in \mathbb{N} \cup \{+\infty\}$.

Definition 47. Underlying structure

An (*underlying*) *structure* of a signal s over a set V , is a graph G with vertex set V .

Remark. Example of images, time series and graph signals

When there is a unique clear underlying structure, we say it is *the* underlying structure. For example, images are compactly supported signals over \mathbb{Z}^2 and their underlying structure is a rectangular grid graph. Time series are signals over \mathbb{N} and their underlying structure is a line graph. The underlying structure of a graph signal is obviously the graph itself.

1.3.2 Graphs in deep learning

TODO: below

We come across the notion of graphs several times in deep learning:

- Connections between two layers of a deep learning model can be represented as a bipartite graph, the *connectivity graph*. It encodes how

the information is propagated through a layer to another. See Section 1.3.2.1.

- Neural architectures can be represented by a graph. In particular, a computation graph is used by deep learning programming languages to keep track of the dependencies between layers of a deep learning model, in order to compute forward and back-propagation. See Section 1.3.2.2.
- A graph can represent the underlying structure of an object (often a vector or a signal). The nodes represent its features, and the edges represent some structural property. See Section 1.3.2.3.
- Datasets can also be graph-structured. The nodes represent the objects of the dataset, and its edge represent some sort of relation between them. See Section 1.3.2.4.

1.3.2.1 Connectivity graph

A Connectivity graph is the bipartite graph whose adjacency matrix is the connectivity matrix of a layer of neurons. Formally, given a linear part of a layer, let \mathbf{x} and \mathbf{y} be the input and output signals, n the size of the set of input neurons $N = \{u_1, u_2, \dots, u_n\}$, and m the size of the set of output neurons $M = \{v_1, v_2, \dots, v_m\}$. This layer implements the equation $y = \Theta x$ where Θ is a $n \times m$ matrix.

Definition 48. The *connectivity graph* $G = (V, E)$ is defined such that $V = N \cup M$ and $E = \{(u_i, v_j) \in N \times M, \Theta_{ij} \neq 0\}$.

I.e. the connectivity graph is obtained by drawing an edge between neurons for which $\Theta_{ij} \neq 0$. For instance, in the special case of a complete bipartite graph, we would obtain a dense layer. Connectivity graphs are especially useful to represent partially connected layers, for which most of the Θ_{ij} are 0. For example, in the case of layers characterized by a small local receptive field, the connectivity graph would be sparse, and output neurons would be

connected to a set of input neurons that corresponds to features that are close together in the input space. Figure 1.3 depicts some examples.

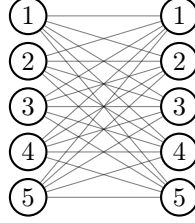


Figure 1.3: Examples

TODO: Figure 1.3. It's just a placeholder right now

Connectivity graphs also allow to graphically modelize how weights are tied in a neural layer. Let's suppose the $\Theta_{i,j}$ are taking their values only into the finite set $K = \{w_1, w_2, \dots, w_\kappa\}$ of size κ , which we will refer to as the *kernel of weights*. Then we can define a labelling of the edges $s : E \rightarrow K$. s is called the *weight sharing scheme* of the layer. This layer can then be formulated as $\forall v \in M, y_v = \sum_{u \in N, (u,v) \in E} w_{s(u,v)} x_u$. Figure 1.4 depicts the connectivity graph of a 1-d convolution layer and its weight sharing scheme.

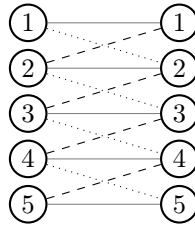


Figure 1.4: Depiction of a 1D-convolutional layer and its weight sharing scheme.

TODO: Add weight sharing scheme in Figure 1.4

1.3.2.2 Computation graph

1.3.2.3 Underlying graph structure and signals

1.3.2.4 Graph-structured dataset

Chapter 2

Convolution on graph domains

Defining a convolution of signals over graph domains is a challenging problem. Obviously, if the graph is not a grid graph there exists no natural definition. We first analyse the reasons why the classical convolution operator is useful in deep learning, and characterize them. Then we will search for domains onto which a convolution with this characteristics can be naturally obtained. This will lead us to put our interest on representation theory and convolutions defined on algebraic structures, such as groups, in order to transfer its construction on vertex domains of graphs, agnostically of the edge set. Then, we will introduce the role of the edge sets and how it should influence the construction. This will provide us with some particular classes of graphs for which we will obtain a natural construction with respect to the aforementioned wanted characteristics. We finally study how we can loosen them to adapt the construction to more irregular domains.

2.1	Analysis of the classical convolution	40
2.2	Construction on the vertex set through groups . .	45
2.3	Construction with the edge set	53
2.4	Conclusion of chapter 2	59

2.1 Analysis of the classical convolution

In this section, we are exposing a few properties of the classical convolution that a generalization to graphs would likely try to preserve. For now let's consider a graph G agnostically of its edges *i.e.* $G \cong V$ is just the set of its vertices.

2.1.1 Properties of the convolution

Consider an edge-less grid graph *i.e.* $G \cong \mathbb{Z}^2$. By restriction to compactly supported signals, this case encompass the case of images.

Definition 49. Convolution on $\mathcal{S}(\mathbb{Z}^2)$

Recall that the convolution between two signals s_1 and s_2 over \mathbb{Z}^2 is a binary operator in $\mathcal{S}(\mathbb{Z}^2)$ defined as:

$$\forall (a, b) \in \mathbb{Z}^2, (s_1 * s_2)[a, b] = \sum_i \sum_j s_1[i, j] s_2[a - i, b - j]$$

TODO: list some properties

2.1.2 Characterization on grid graphs

Let's recall first what is a transformation, and how it acts on signals.

Definition 50. Transformation

A *transformation* $f : V \rightarrow V$ is a function with same domain and codomain. The set of transformations is denoted $\Phi(V)$. The set of bijective transformations is denoted $\Phi^*(V) \subset \Phi(V)$.

In case $f \in \Phi^*(V)$, it can act on $\mathcal{S}(V)$ through the linear operator $L_f \in \mathcal{L}(\mathcal{S}(V))$ defined as:

$$\forall s \in \mathcal{S}(V), \forall v \in V, f(s)[v] := L_f s[v] = s[f^{-1}(v)]$$

i.e. an entry of a transformed signal is obtained by doing a lookup of the entry of the original signal.

In case $f \notin \Phi^*(V)$, we can still define $L_f \in \mathcal{L}(\mathcal{S}(V))$, however we need to linearly aggregate the entries on the fibers:

$$\forall s \in \mathcal{S}(V), \forall v \in V, f(s)[v] := L_f s[v] = \text{AGGREGATE}\{s[u], u \in f^{-1}\{v\}\}$$

where AGGREGATE can be for example the sum, the average, or the max, and $\text{AGGREGATE}(\emptyset) = 0$.

We also recall the formalism of translations.

Definition 51. Translation on $\mathcal{S}(\mathbb{Z}^2)$

A translation on \mathbb{Z}^2 is defined as a transformation $t \in \Phi^*(\mathbb{Z}^2)$ such that

$$\exists(a, b) \in \mathbb{Z}^2, \forall(x, y) \in \mathbb{Z}^2, t(x, y) = (x + a, y + b)$$

It also acts on $\mathcal{S}(\mathbb{Z}^2)$ with the notation $t_{a,b}$ *i.e.*

$$\forall s \in \mathcal{S}(\mathbb{Z}^2), \forall(x, y) \in \mathbb{Z}^2, t_{a,b}(s)[x, y] = s[x - a, y - b]$$

For any set E , we denote by $\mathcal{T}(E)$ its translations if they are defined.

The next proposition fully characterizes convolution operators with their translational equivariance property. This can be seen as a discretization of a classic result from the theory of distributions.

Proposition 52. Characterization of convolution operators on $\mathcal{S}(\mathbb{Z}^2)$

On real-valued signals over \mathbb{Z}^2 , the class of linear transformations that are equivariant to translations is exactly the class of convolutive operations *i.e.*

$$\exists w \in \mathcal{S}(\mathbb{Z}^2), f = . * w \Leftrightarrow \begin{cases} f \in \mathcal{L}(\mathcal{S}(\mathbb{Z}^2)) \\ \forall t \in \mathcal{T}(\mathcal{S}(\mathbb{Z}^2)), f \circ t = t \circ f \end{cases}$$

Proof. The result from left to right is a direct consequence of the definitions:

$$\begin{aligned}
& \forall s \in \mathcal{S}(\mathbb{Z}^2), \forall s' \in \mathcal{S}(\mathbb{Z}^2), \forall (\alpha, \beta) \in \mathbb{R}^2, \forall (a, b) \in \mathbb{Z}^2, \\
& f_w(\alpha s + \beta s')[a, b] = \sum_i \sum_j (\alpha s + \beta s')[i, j] w[a - i, b - j] \\
& \quad = \alpha f_w(s)[a, b] + \beta f_w(s')[a, b] \quad (\text{linearity}) \\
& \forall s \in \mathcal{S}(\mathbb{Z}^2), \forall (\alpha, \beta) \in \mathbb{Z}^2, \forall (a, b) \in \mathbb{Z}^2, \\
& f_w \circ t_{\alpha, \beta}(s)[a, b] = \sum_i \sum_j t_{\alpha, \beta}(s)[i, j] w[a - i, b - j] \\
& \quad = \sum_i \sum_j s[i - \alpha, j - \beta] w[a - i, b - j] \\
& \quad = \sum_{i'} \sum_{j'} s[i', j'] w[a - i' - \alpha, b - j' - \beta] \quad (2.1) \\
& \quad = f_w(s)[a - \alpha, b - \beta] \\
& \quad = t_{\alpha, \beta} \circ f_w(s)[a, b] \quad (\text{equivariance})
\end{aligned}$$

Now let's prove the result from right to left .

Let $f \in \mathcal{L}(\mathcal{S}(\mathbb{Z}^2))$, $s \in \mathcal{S}(\mathbb{Z}^2)$. We suppose that f commutes with translations. Recall that s can be linearly decomposed on the infinite family of dirac signals:

$$s = \sum_i \sum_j s[i, j] \delta_{i, j}, \text{ where } \delta_{i, j}[x, y] = \begin{cases} 1 & \text{if } (x, y) = (i, j) \\ 0 & \text{otherwise} \end{cases}$$

By linearity of f and then equivariance to translations:

$$\begin{aligned}
f(s) &= \sum_i \sum_j s[i, j] f(\delta_{i, j}) \\
&= \sum_i \sum_j s[i, j] f \circ t_{i, j}(\delta_{0, 0}) \\
&= \sum_i \sum_j s[i, j] t_{i, j} \circ f(\delta_{0, 0})
\end{aligned}$$

By denoting $w = f(\delta_{0,0}) \in \mathcal{S}(\mathbb{Z}^2)$, we obtain:

$$\begin{aligned} \forall (a, b) \in \mathbb{Z}^2, f(s)[a, b] &= \sum_i \sum_j s[i, j] t_{i,j}(w)[a, b] \\ &= \sum_i \sum_j s[i, j] w[a - i, b - j] \\ \text{i.e. } f(s) &= s * w \end{aligned} \tag{2.2}$$

□

2.1.3 Usefulness of convolutions in deep learning

Remark. **Equivariance property of CNNs**

In deep learning, an important argument in favor of CNNs is that convolutional layers are equivariant to translations. Intuitively, that means that a detail of an object in an image should produce the same features independently of its position in the image.

Remark. **Lossless superiority of CNNs over MLPs**

The converse result, as a consequence of Proposition 75, is never mentioned in deep learning literature. However it is also a strong one: it means that layers of CNNs have every translational equivariant functions in their search space, so it implies that the reduction of parameters from an MLP to a CNN is done with strictly no loss of expressivity (provided the objective function is known to bear this property). Besides, it helps the training to search in a much more confined space.

Hence, in our construction, we will try to preserve the characterization from Proposition 75 as it is mostly the reason why they are successful in deep learning. Note that the reduction of parameters compared to a dense layer is also a consequence of this characterization.

2.1.4 Related work generalising the discrete convolution on graphs

TODO:

2.2 Construction on the vertex set through groups

As Proposition 75 is a complete characterization of convolutions, it can be used to define them *i.e.* convolutive operations can be constructed as the set of linear transformations that are equivariant to translations. However, in the general case where G is not a grid graph, translations are not defined, so that construction needs to be generalized beyond translational equivariances.

TODO: reword sentence below (not exactly true)

A classic result from group theory is that one sense of this characterization (equivariance) hold for groups of (bijective) transformations, where the group of translations is a particular case. More generally, it also holds for groupoids (Weinstein, 1996). The converse sense will be discussed later.

Note that our approach is different than in (Cohen and Welling, 2016), where the authors define group-equivariant convolutions from an already defined convolution. On graphs, convolutions aren't already defined. We are studying group convolutions in view of constructing one.

Definition 53. Group convolution I

Let a group \mathcal{G} , the group convolution I between two signals s_1 and $s_2 \in \mathcal{S}(\mathcal{G})$ is defined as:

$$\forall h \in \mathcal{G}, (s_1 *_I s_2)[h] = \sum_{g \in \mathcal{G}} s_1[g] s_2[g^{-1}h]$$

provided one of the signals has finite support if \mathcal{G} is not finite.

For a graph $G = \langle V, E \rangle$ and a subgroup $\Gamma \subset \Phi^*(V)$, this definition is applicable for $\mathcal{S}(\Gamma)$, but not for $\mathcal{S}(V)$ as V is not a group. Nonetheless, our point here is that we will use the group convolution on $\mathcal{S}(\Gamma)$ to construct the convolutions on $\mathcal{S}(V)$.

For now, let's assume Γ is in one-to-one correspondence with V , and let's define an isomorphism φ from Γ to V , to equip V with a group structure. We denote $\Gamma \xrightarrow{\varphi} V$ and $g_v \mapsto v$.

Then, the linear morphism $\tilde{\varphi}$ from $\mathcal{S}(\Gamma)$ to $\mathcal{S}(V)$ defined on the Dirac bases by $\tilde{\varphi}(\delta_g) = \delta_{\varphi(g)}$ is also an isomorphism. Hence, V and $\mathcal{S}(V)$ would inherit the same inherent structural properties as Γ and $\mathcal{S}(\Gamma)$, in addition of being related in the same way. For notational simplicity, we will use the same symbol φ for φ and $\tilde{\varphi}$ (as done between f and L_f). A commutative diagram between the sets is depicted on Figure 2.1.

$$\begin{array}{ccc} \Gamma & \xrightarrow{\varphi} & V \\ s \downarrow & & \downarrow s \\ \mathcal{S}(\Gamma) & \xrightarrow{\varphi} & \mathcal{S}(V) \end{array}$$

Figure 2.1: Commutative diagram between sets

Lemma 54. Relation between $\mathcal{S}(\Gamma)$ and $\mathcal{S}(V)$

$$\forall s \in \mathcal{S}(\Gamma), \forall u \in V, \varphi(s)[u] = s[\varphi^{-1}(u)] = s[g_u]$$

Proof.

$$\begin{aligned} \forall s \in \mathcal{S}(\Gamma), \varphi(s) &= \varphi\left(\sum_{g \in \Gamma} s[g] \delta_g\right) = \sum_{g \in \Gamma} s[g] \varphi(\delta_g) = \sum_{g \in \Gamma} s[g] \delta_{\varphi(g)} \\ &= \sum_{v \in V} s[g_v] \delta_v \end{aligned}$$

$$\text{So } \forall v \in V, \varphi(s)[u] = s[g_u]$$

□

Hence, we can steer the definition of the group convolution from $\mathcal{S}(\Gamma)$ to $\mathcal{S}(V)$ as follows:

Definition 55. Group convolution II

Let a subgroup $\Gamma \subset \Phi^*(V)$ such that $\Gamma \cong^{\varphi} V$. The group convolution II

2.2. CONSTRUCTION ON THE VERTEX SET THROUGH GROUPS 47

between two signals s_1 and $s_2 \in \mathcal{S}(V)$ is defined as:

$$\begin{aligned} \forall u \in V, (s_1 *_{\text{II}} s_2)[u] &= \sum_{v \in V} s_1[v] s_2[\varphi(g_v^{-1} g_u)] \\ &= \sum_{\substack{(a,b) \in V^2 \\ \text{s.t. } g_a g_b = g_u}} s_1[a] s_2[b] \end{aligned}$$

Proposition 56. Relation between group convolution I and II

Let a subgroup $\Gamma \subset \Phi^*(V)$ such that $\Gamma \stackrel{\varphi}{\cong} V$,

$$\forall s_1, s_2 \in \mathcal{S}(\Gamma), \forall u \in V, (\varphi(s_1) *_{\text{II}} \varphi(s_2))[u] = (s_1 *_I s_2)[g_u]$$

Proof. Using Lemma 54,

$$\begin{aligned} (\varphi(s_1) *_{\text{II}} \varphi(s_2))[u] &= \sum_{v \in V} \varphi(s_1)[v] \varphi(s_2)[\varphi(g_v^{-1} g_u)] \\ &= \sum_{v \in V} s_1[g_v] s_2[g_v^{-1} g_u] \\ &= \sum_{g \in \Gamma} s_1[g] s_2[g^{-1} g_u] \\ &= (s_1 *_I s_2)[g_u] \end{aligned}$$

□

Proposition 57. Equivariance to $\varphi(\Gamma)$

With Definition 55, convolution operators acting on the right of $\mathcal{S}(V)$ are equivariant to $\varphi(\Gamma)$ i.e.

$$\exists w \in \mathcal{S}(V), f = . *_{\text{II}} w \Rightarrow \forall v \in V, f \circ \varphi(g_v) = \varphi(g_v) \circ f$$

Proof.

$$\begin{aligned}
\forall s \in \mathcal{S}(V), \forall u \in V, \forall v \in V, \\
(f_w \circ \varphi(g_u))(s)[v] &= \sum_{\substack{(a,b) \in V^2 \\ s.t. \ g_a g_b = g_v}} \varphi(g_u)(s)[a] w[b] \\
&= \sum_{\substack{(a,b) \in V^2 \\ s.t. \ g_a g_b = g_v}} s[\varphi(g_u)^{-1}(a)] w[b] \\
&= \sum_{\substack{(a,b) \in V^2 \\ s.t. \ g_{\varphi(g_u)(a)} g_b = g_v}} s[a] w[b]
\end{aligned}$$

Because φ is an isomorphism, its inverse $c \mapsto g_c$ is also an isomorphism and so $g_{\varphi(g_u)(a)} g_b = g_v \Leftrightarrow g_a g_b = g_{\varphi(g_u)^{-1}(v)}$. So we have:

$$\begin{aligned}
(f_w \circ \varphi(g_u))(s)[v] &= \sum_{\substack{(a,b) \in V^2 \\ s.t. \ g_a g_b = g_{\varphi(g_u)^{-1}(v)}}} s[a] w[b] \\
&= s *_I w[\varphi(g_u)^{-1}(v)] \\
&= (\varphi(g_u) \circ f_w)(s)[v]
\end{aligned}$$

□

Remark. Note that convolution operators of the form $f_w = . *_I w$ are also equivariant to \mathcal{G} , but the proposition and the proof are omitted as they are similar to the latter.

In fact, both group convolutions are the same as the latter one borrows the algebraic structure of the first one. Thus we only obtain equivariance to $\varphi(\Gamma)$, which is actually V equipped with the group structure of Γ mirrored via φ . To obtain equivariance to Γ , we need to take into account the fact that it contains bijective transformations of V .

2.2. CONSTRUCTION ON THE VERTEX SET THROUGH GROUPS 49

Hence, note that $g \in \Gamma$ can act on both Γ through the left multiplication and on V as being an object of $\Phi^*(V)$. This ambivalence can be seen on a commutative diagram, see Figure 2.2.

$$\begin{array}{ccc} g_u & \xrightarrow{g_v} & g_v g_u \\ \varphi \downarrow & & \downarrow \varphi \\ u & \xrightarrow[g_v]{(P)} & \varphi(g_v g_u) \end{array}$$

Figure 2.2: Commutative diagram. All arrows except for the one labeled with (P) are always True.

For (P) to be true means that φ is an equivariant map between groups *i.e.* whether the mapping is done before or after the action of Γ has no impact on the result. When such φ exists, Γ and V are said to be equivalent (in the isomorphic sense) and we denote $\Gamma \equiv V$.

Definition 58. Equivariant map

An isomorphism φ between a group of permutations \mathcal{G} and a group \mathcal{H} is said to be an *equivariant map* if

$$\forall g, h \in \mathcal{G}, g(\varphi(h)) = \varphi(gh)$$

In our case, if $\Gamma \equiv V$, we are interested to know if our φ exhibits the equivalence.

Definition 59. φ -Equivalence

A subgroup $\Gamma \subset \Phi^*(V)$ such that $\Gamma \stackrel{\varphi}{\cong} V$, is said to be *φ -equivalent* if φ is an equivariant map *i.e.* if it verifies the property:

$$\forall v, u \in V, g_v(u) = \varphi(g_v g_u) \tag{P}$$

In that case we denote $\Gamma \stackrel{\varphi}{\equiv} V$

Remark. For example, translations on the grid graph, with $\varphi(t_{i,j}) = (i, j)$, are φ -equivalent as $t_{i,j}(a, b) = \varphi(t_{i,j} \circ t_{a,b})$. However, with $\varphi(t_{i,j}) = (-i, -j)$, they would not be φ -equivalent.

Definition 60. Group convolution III

Let a subgroup $\Gamma \subset \Phi^*(V)$ such that $\Gamma \stackrel{\varphi}{\cong} V$. The group convolution III between two signals s_1 and $s_2 \in \mathcal{S}(V)$ is defined as:

$$s_1 *_{\text{III}} s_2 = \sum_{v \in V} s_1[v] g_v(s_2) \quad (2.3)$$

$$= \sum_{g \in \Gamma} s_1[\varphi(g)] g(s_2) \quad (2.4)$$

The two expressions differ on the domain upon which the summation is done. The expression (2.3) put the emphasis on each vertex and its action, whereas the expression (2.4) emphasizes on each object of Γ .

Lemma 61. Relation with group convolution II

$$\Gamma \stackrel{\varphi}{\cong} V \Leftrightarrow *_{\text{II}} = *_{\text{III}}$$

Proof.

$$\begin{aligned} \forall s_1, s_2 \in \mathcal{S}(V), \\ s_1 *_{\text{II}} s_2 &= s_1 *_{\text{III}} s_2 \\ \Leftrightarrow \forall u \in V, \sum_{v \in V} s_1[v] s_2[\varphi(g_v^{-1} g_u)] &= \sum_{v \in V} s_1[v] s_2[g_v^{-1}(u)] \end{aligned} \quad (2.5)$$

Hence, the direct sense is obtained by applying (P).

For the converse, given $u, v \in V$, we first realize (2.5) for $s_1 := \delta_v$, obtaining $s_2[\varphi(g_v^{-1} g_u)] = s_2[g_v^{-1}(u)]$, which we then realize for a real signal s_2 having no two equal entries, obtaining $\varphi(g_v^{-1} g_u) = g_v^{-1}(u)$. From the latter we finally obtain (P) with the one-to-one correspondence $v := v^{-1}$, where $v^{-1} = \varphi(g_v^{-1})$ and using the fact that φ and φ^{-1} are isomorphisms. \square

2.2. CONSTRUCTION ON THE VERTEX SET THROUGH GROUPS 51

This time, we obtain equivariance to Γ .

Proposition 62. Equivariance to Γ

If Γ is φ -equivalent, convolution operators acting on the right of $\mathcal{S}(V)$ are equivariant to Γ *i.e.*

$$\begin{cases} \Gamma \stackrel{\varphi}{\equiv} V \\ \exists w \in \mathcal{S}(V), f = . * w \end{cases} \Rightarrow \forall g \in \Gamma, f \circ g = g \circ f$$

where $*$ denote either $*_{\text{II}}$ or $*_{\text{III}}$.

Proof. In the following equations, (2.6) is obtained from Lemma 61, (2.7) is obtained because left multiplication in a group is an isomorphism, and (2.8) is obtained because of (P).

$$\begin{aligned} \forall g \in \Gamma, \forall s \in \mathcal{S}(V), \\ f_w \circ g(s) &= \sum_{h \in \Gamma} g(s)[\varphi(h)] h(w) \end{aligned} \quad (2.6)$$

$$= \sum_{h \in \Gamma} g(s)[\varphi(gh)] gh(w) \quad (2.7)$$

$$\begin{aligned} &= \sum_{h \in \Gamma} g(s)[g(\varphi(h))] gh(w) \quad (2.8) \\ &= \sum_{h \in \Gamma} s[\varphi(h)] gh(w) \end{aligned}$$

$$\begin{aligned} \forall g \in \Gamma, \forall s \in \mathcal{S}(V), \forall u \in V, \\ f_w \circ g(s)[u] &= \sum_{h \in \Gamma} s[\varphi(h)] h(w)[g^{-1}(u)] \\ &= f_w(s)[g^{-1}(u)] \\ &= g \circ f_w(s)[u] \end{aligned}$$

□

From (2.4), we define a mixed domain convolution *i.e.* that is defined for

$r \in \mathcal{S}(\Gamma)$ and $s \in \mathcal{S}(V)$, without the need of expliciting the isomorphisms φ .

Definition 63. Mixed domain convolution

Let a subgroup $\Gamma \subset \Phi^*(V)$ such that $V \cong \Gamma$. The mixed domain convolution between two signals $r \in \mathcal{S}(\Gamma)$ and $s \in \mathcal{S}(V)$ results in a signal $r *_{\text{M}} s \in \mathcal{S}(V)$ and is defined as:

$$\forall u \in V, (r *_{\text{M}} s)[u] = \sum_{g \in \Gamma} r[g] g(s)$$

Proposition 64. Relation with group convolution III

$$\forall \varphi \in \text{ISO}(\Gamma, V), r *_{\text{M}} s = \varphi(r) *_{\text{H}} s$$

Proof. **TODO:**

□

Proposition 65. Relation with group convolution I

Let $\varphi \in \text{ISO}(\Gamma, V)$, we have

$$r *_{\text{M}} s = r *_{\text{I}} \varphi^{-1}(s) \Leftrightarrow \Gamma \stackrel{\varphi}{\cong} V$$

Proof. **TODO:**

□

The direct result is important because it justifies that even without the need of expliciting them, there are underlying isomorphisms φ for which the group Γ is φ -equivalent.

2.3 Construction with the edge set

At this point, there are a two drawbacks: (D1) that $V \cong \Gamma$, (D2) the converse of the characterization doesn't hold. Let's now consider the edge set in our construction and keep these drawbacks in mind.

2.3.1 Construction on Cayley graphs

We now consider the class of Cayley (di)graphs (Cayley, 1878; Wikipedia, 2018a) because the φ -coherence property (P) naturally holds on them.

TODO: rewrite and obtain: On a Cayley graph, there is φ such that the subgroupoid of transformation that are EC is φ -coherent (proof on \mathcal{U} then extended)

Definition 66. Cayley graph

Let a group Γ and one of its generating set \mathcal{U} . The *Cayley graph* generated by \mathcal{U} , is the digraph $\vec{G} = \langle V, E \rangle$ such that $V \cong \Gamma$ through $u \xrightarrow{\varphi^{-1}} g_u$, and E is such that:

$$a \sim b \Leftrightarrow \exists u \in \varphi(\mathcal{U}) \subset V, g_u g_a = g_b$$

Proposition 67. A Cayley graph of group $\Gamma \xrightarrow{\varphi} V$ is φ -coherent.

Proof. Let $\mathcal{U} \subset \Gamma$ the generating set of the Cayley graph.

TODO:

□

We can alleviate (D1) by summing onto the generating set \mathcal{U} instead of onto Γ , which also makes the convolution on Cayley graphs edge-constrained.

TODO: not equivariant unless abelian group

Definition 68. Cayley graph convolution

$$\forall u \in V, (s_1 * s_2)[u] = \sum_{g \in \mathcal{U}} s_1[g] g(s_2)$$

Conversely, let a graph $G = \langle V, E \rangle$ and a φ -coherent subgroup $\Gamma \subseteq \Phi^*(V)$. Then we can generate a Cayley graphs with a generating set of Γ , and thus define a Cayley graph convolution. In particular, when transformations of Γ are edge-constrained, that would also be the case for this convolution on G , which leads us to study edge-constrained transformations.

TODO: operator and characterization

TODO: which graph is a Cayley graph ?

2.3.2 Construction on graph groupoids

TODO: work in progress

On graphs, we notice that the property (P) can be realized by transformations acting on edges. However, unless the graph is complete, these actions can't be composed everywhere to form another edge constrained action. The algebraic structure that possesses the same kind of properties than a group except that its composition law is not defined everywhere is called a groupoid. The following definitions clarify our discussion.

Definition 69. Groupoid

A groupoid is a set equipped with a closed partial composition law, a unique identity element, and every unique inverses.

Remark. We use the convention than left and right inverses must be the same.

Definition 70. Graph groupoid

The *groupoid* $\mathcal{P}(G)$ of a graph $G = \langle V, E \rangle$ is the set of its paths equipped with:

1. two maps ψ and φ that respectively map a path to its first and last element,
2. a closed partial composition law gh defined if and only if $\psi(g) = \varphi(h)$, which concatenates g behind h and ~~removes adjacent duplicate vertices~~ to rewrite,
3. an inverse operator $^{-1}$ which maps a path to its reverse,
4. an identity element Id which is the path of length 0.

Remark. Recall from Definition 34 that a path can't contain adjacent duplicates.

Remark. Note that even though the composite path gh has elements of h before those of g we write gh instead of hg because we'll need the left operand to act on the right one through functional notation $g(h)$.

Definition 71. Graph k -groupoid

The *k -groupoid* $\mathcal{P}_k(G)$ of a graph $G = \langle V, E \rangle$, for $k \in \mathbb{N}^*$, is the groupoid obtained by restricting $\mathcal{P}(G)$ to paths of length at most k (the definition domain of its composition law is also further restricted by the length of the resulting paths in $\mathcal{P}(G)$).

Definition 72. k -Groupoid convolution

Let a graph $G = \langle V, E \rangle$. Let a subgroupoid $\Gamma \subseteq \mathcal{P}_k(G)$. The *k -groupoid*

convolution between two signals s_1 and $s_2 \in \mathcal{S}(\Gamma)$ is defined as:

$$\begin{aligned}
 \forall h \in \Gamma, (s_1 * s_2)[h] &= \sum_{\substack{(a,b) \in \Gamma^2 \\ s.t. \ ab=h}} s_1[a] s_2[b] \\
 &= \sum_{\substack{g \in \Gamma \\ s.t. \ \varphi(g)=\varphi(h)}} s_1[g] s_2[g^{-1}h] \\
 &= \sum_{\substack{g \in \Gamma \\ s.t. \ \psi(g)=\psi(h)}} s_1[hg^{-1}] s_2[g]
 \end{aligned}$$

Claim 73. Path transformation

Let a graph $G = \langle V, E \rangle$. By identifying vertices with paths of length 1, a path $g \in \mathcal{P}(G)$ can act as a transformation on $v \in V$ through the composition law of $\mathcal{P}(G)$. Also note that $g(v) = g(v^{-1})$.

We can now define the k -Groupoid convolution operator on $\mathcal{S}(G)$ by restriction of the second operand from $\mathcal{S}(\Gamma)$ to paths of length 1:

Definition 74. k -Groupoid convolution operator

Let a graph $G = \langle V, E \rangle$. Let a subgroupoid $\Gamma \subseteq \mathcal{P}_k(G)$. The k -groupoid convolution operator f_w with parameter $w \in \mathcal{S}(\Gamma)$ is defined as:

$$\forall s \in \mathcal{S}(\Gamma), \forall h \in \Gamma, f_w(s)[h] = (s * w)[h]$$

And when restricted to $\mathcal{S}(G)$ it is defined as:

$$\begin{aligned}
 \forall s \in \mathcal{S}(G), \forall v \in V, f_w(s)[v] &= \sum_{\substack{g \in \Gamma \\ s.t. \ \psi(g)=v}} s[g(v)] w[g] \\
 \forall s \in \mathcal{S}(G), \forall v \in V, f_w(s)[v] &= \sum_{\substack{g \in \Gamma \\ s.t. \ \varphi(g)=v}} s[g] w[g^{-1}(v)]
 \end{aligned}$$

Proposition 75. Groupoid equivariance to Γ

k -Groupoid convolution operators on $\mathcal{S}(G)$ are groupoid equivariant to Γ i.e.

$$\exists w \in \mathcal{S}(\Gamma), f = w * . \Rightarrow \forall v \in V, \forall g \in \Gamma \text{ s.t. } \psi(g^{-1}) = v, f \circ g[v] = g \circ f[v]$$

$$g(h(v)) \text{ maybe false} \tag{2.9}$$

Mini patron of todo:

- Equivariance to Γ holds, proof
- Converse of characterization does not hold yet, except on orbits
- property for it to hold
- relaxing one-to-one correspondence constraint but keeping other properties
- other avenue instead of property: should make use of edges to build a group structure
- ideal graph (lattice-regular)
- if group is too much then just groupoid structure from edges is enough

TODO: finish this section

2.3.3 To rename**Definition 76. Graph automorphisms**

A graph automorphism of a graph $G = \langle V, E \rangle$ is a bijection in the vertex

domain $\phi : V \rightarrow V$ such that $\{u, v\} \in E \Leftrightarrow \{\phi(u), \phi(v)\} \in E$. We denote $\mathcal{A}(G)$ the group of automorphism on G .

We denote by $\mathcal{E}(\phi)$ the set of input-output mapping of ϕ , defined as $\mathcal{E}(\phi) = \{(x, y) \in V^2, \phi(x) = y\}$.

A graph automorphism ϕ is said to be *edge-constrained* (EC) if $\mathcal{E}(\phi) \subseteq E$. We denote $\mathcal{A}_{\text{ec}}(G)$ the set of edge-constrained automorphism on G .

Definition 77. Orthogonality

Two graph automorphisms ϕ_1 and ϕ_2 are said to be orthogonal, if and only if $\mathcal{E}(\phi_1) \cap \mathcal{E}(\phi_2) = \emptyset$, denoted $\phi_1 \perp \phi_2$. They are said to be aligned otherwise. Similarly, we define orthogonality of r automorphisms as $\phi_1 \perp \cdots \perp \phi_r \Leftrightarrow \mathcal{E}(\phi_1) \cap \cdots \cap \mathcal{E}(\phi_r) = \emptyset$

2.3.4 Lattice-regular graph

Definition 78. Lattice-regular graph

A lattice-regular graph is a regular graph that admits r orthogonal edge-constrained automorphisms, where r is its degree.

2.4 Conclusion of chapter 2

TODO:

Chapter 3

Neural networks on graphs

Chapter 4

Industrial applications

TODO:

Conclusion

TODO:

Bibliography

- Abadi, Martín, Ashish Agarwal, Paul Barham, Eugene Brevdo, Zhifeng Chen, Craig Citro, Greg S. Corrado, Andy Davis, Jeffrey Dean, Matthieu Devin, Sanjay Ghemawat, Ian Goodfellow, Andrew Harp, Geoffrey Irving, Michael Isard, Yangqing Jia, Rafal Jozefowicz, Lukasz Kaiser, Manjunath Kudlur, Josh Levenberg, Dan Mané, Rajat Monga, Sherry Moore, Derek Murray, Chris Olah, Mike Schuster, Jonathon Shlens, Benoit Steiner, Ilya Sutskever, Kunal Talwar, Paul Tucker, Vincent Vanhoucke, Vijay Vasudevan, Fernanda Viégas, Oriol Vinyals, Pete Warden, Martin Wattenberg, Martin Wicke, Yuan Yu, and Xiaoqiang Zheng (2015). *TensorFlow: Large-Scale Machine Learning on Heterogeneous Systems*. Software available from tensorflow.org. URL: <http://tensorflow.org/> (cit. on p. 25).
- Arora, Raman, Amitabh Basu, Poorya Mianjy, and Anirbit Mukherjee (2018). “Understanding Deep Neural Networks with Rectified Linear Units”. In: *International Conference on Learning Representations*. URL: https://openreview.net/forum?id=B1J_rgWRW (cit. on pp. 17, 18).
- Bass, Jean (1968). “Cours de mathématiques”. In: (cit. on p. 8).
- Bengio, Yoshua (2009). “Learning deep architectures for AI”. In: *Foundations and trends® in Machine Learning* 2.1, pp. 1–127 (cit. on p. 17).
- Bianchini, Monica and Franco Scarselli (2014). “On the complexity of neural network classifiers: A comparison between shallow and deep architectures”. In: *IEEE transactions on neural networks and learning systems* 25.8, pp. 1553–1565 (cit. on p. 18).

- Bottou, Léon (2010). “Large-scale machine learning with stochastic gradient descent”. In: *Proceedings of COMPSTAT’2010*. Springer, pp. 177–186 (cit. on p. 30).
- Cayley, Professor (1878). “Desiderata and Suggestions: No. 2. The Theory of Groups: Graphical Representation”. In: *American Journal of Mathematics* 1.2, pp. 174–176. ISSN: 00029327, 10806377. URL: <http://www.jstor.org/stable/2369306> (cit. on p. 53).
- Clevert, Djork-Arné, Thomas Unterthiner, and Sepp Hochreiter (2015). “Fast and accurate deep network learning by exponential linear units (elus)”. In: *arXiv preprint arXiv:1511.07289* (cit. on p. 16).
- Cohen, Nadav and Amnon Shashua (2016). “Convolutional rectifier networks as generalized tensor decompositions”. In: *International Conference on Machine Learning*, pp. 955–963 (cit. on p. 18).
- Cohen, Nadav, Ronen Tamari, and Amnon Shashua (2018a). “Boosting Dilated Convolutional Networks with Mixed Tensor Decompositions”. In: *International Conference on Learning Representations*. URL: <https://openreview.net/forum?id=S1JHhv6TW> (cit. on pp. 17, 19).
- Cohen, Taco and Max Welling (2016). “Group equivariant convolutional networks”. In: *International Conference on Machine Learning*, pp. 2990–2999 (cit. on pp. 45, 84).
- Cohen, Taco S., Mario Geiger, Jonas Köhler, and Max Welling (2018b). “Spherical CNNs”. In: *International Conference on Learning Representations*. URL: <https://openreview.net/forum?id=Hkbd5xZRb> (cit. on p. 84).
- Cybenko, George (1989). “Approximation by superpositions of a sigmoidal function”. In: *Mathematics of control, signals and systems* 2.4, pp. 303–314 (cit. on p. 16).
- Delalleau, Olivier and Yoshua Bengio (2011). “Shallow vs. deep sum-product networks”. In: *Advances in Neural Information Processing Systems*, pp. 666–674 (cit. on pp. 17, 18).

- Deng, Jia, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei (2009). “Imagenet: A large-scale hierarchical image database”. In: *Computer Vision and Pattern Recognition, 2009. CVPR 2009. IEEE Conference on*. IEEE, pp. 248–255 (cit. on pp. 17, 19).
- Eldan, Ronen and Ohad Shamir (2016). “The power of depth for feedforward neural networks”. In: *Conference on Learning Theory*, pp. 907–940 (cit. on p. 18).
- Glorot, Xavier and Yoshua Bengio (2010). “Understanding the difficulty of training deep feedforward neural networks”. In: *Proceedings of the thirteenth international conference on artificial intelligence and statistics*, pp. 249–256 (cit. on p. 17).
- Glorot, Xavier, Antoine Bordes, and Yoshua Bengio (2011). “Deep sparse rectifier neural networks”. In: *International Conference on Artificial Intelligence and Statistics*, pp. 315–323 (cit. on pp. 16, 17).
- Hackbusch, Wolfgang (2012). *Tensor spaces and numerical tensor calculus*. Vol. 42. Springer Science & Business Media (cit. on pp. 8, 9).
- Håstad, Johan and Mikael Goldmann (1991). “On the power of small-depth threshold circuits”. In: *Computational Complexity* 1.2, pp. 113–129 (cit. on p. 17).
- He, Kaiming, Xiangyu Zhang, Shaoqing Ren, and Jian Sun (2015). “Delving deep into rectifiers: Surpassing human-level performance on imagenet classification”. In: *Proceedings of the IEEE international conference on computer vision*, pp. 1026–1034 (cit. on p. 16).
- (2016). “Deep residual learning for image recognition”. In: *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 770–778 (cit. on pp. 15, 19).
- Hinton, Geoffrey, Li Deng, Dong Yu, George E Dahl, Abdel-rahman Mohamed, Navdeep Jaitly, Andrew Senior, Vincent Vanhoucke, Patrick Nguyen, Tara N Sainath, et al. (2012). “Deep neural networks for acoustic model-

- ing in speech recognition: The shared views of four research groups”. In: *IEEE Signal Processing Magazine* 29.6, pp. 82–97 (cit. on p. 15).
- Hinton, Geoffrey E, Simon Osindero, and Yee-Whye Teh (2006). “A fast learning algorithm for deep belief nets”. In: *Neural computation* 18.7, pp. 1527–1554 (cit. on p. 17).
- Hoogeboom, Emiel, Jorn W.T. Peters, Taco S. Cohen, and Max Welling (2018). “HexaConv”. In: *International Conference on Learning Representations*. URL: <https://openreview.net/forum?id=r1vuQG-CW> (cit. on p. 84).
- Hornik, Kurt (1991). “Approximation capabilities of multilayer feedforward networks”. In: *Neural networks* 4.2, pp. 251–257 (cit. on p. 16).
- Hornik, Kurt, Maxwell Stinchcombe, and Halbert White (1989). “Multilayer feedforward networks are universal approximators”. In: *Neural Networks* 2.5, pp. 359–366 (cit. on pp. 16, 29).
- Huang, Gao, Zhuang Liu, Kilian Q Weinberger, and Laurens van der Maaten (2017). “Densely connected convolutional networks”. In: *Proceedings of the IEEE conference on computer vision and pattern recognition*. Vol. 1. 2, p. 3 (cit. on pp. 15, 19).
- Jarrett, Kevin, Koray Kavukcuoglu, Yann LeCun, et al. (2009). “What is the best multi-stage architecture for object recognition?” In: *Computer Vision, 2009 IEEE 12th International Conference on*. IEEE, pp. 2146–2153 (cit. on p. 16).
- Klambauer, Günter, Thomas Unterthiner, Andreas Mayr, and Sepp Hochreiter (2017). “Self-Normalizing Neural Networks”. In: *Advances in Neural Information Processing Systems 30*. Ed. by I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett. Curran Associates, Inc., pp. 971–980. URL: <http://papers.nips.cc/paper/6698-self-normalizing-neural-networks.pdf> (cit. on p. 16).
- Krizhevsky, Alex, Ilya Sutskever, and Geoffrey E Hinton (2012). “Imagenet classification with deep convolutional neural networks”. In: *Advances in*

- Neural Information Processing Systems*, pp. 1097–1105 (cit. on pp. 15, 17, 28).
- LeCun, Y. (1987). “Modeles connexionnistes de l’apprentissage (connectionist learning models)”. PhD thesis. Université P. et M. Curie (Paris 6) (cit. on p. 20).
- LeCun, Yann, Yoshua Bengio, et al. (1995). “Convolutional networks for images, speech, and time series”. In: *The handbook of brain theory and neural networks* 3361.10, p. 1995 (cit. on pp. 15, 28).
- LeCun, Yann, Bernhard Boser, John S Denker, Donnie Henderson, Richard E Howard, Wayne Hubbard, and Lawrence D Jackel (1989). “Backpropagation applied to handwritten zip code recognition”. In: *Neural computation* 1.4, pp. 541–551 (cit. on pp. 15, 16).
- LeCun, Yann, Léon Bottou, Yoshua Bengio, and Patrick Haffner (1998). “Gradient-based learning applied to document recognition”. In: *Proceedings of the IEEE* 86.11, pp. 2278–2324 (cit. on p. 29).
- Lin, Henry W, Max Tegmark, and David Rolnick (2017). “Why does deep and cheap learning work so well?” In: *Journal of Statistical Physics* 168.6, pp. 1223–1247 (cit. on p. 18).
- Maas, Andrew L, Awni Y Hannun, and Andrew Y Ng (2013). “Rectifier nonlinearities improve neural network acoustic models”. In: *Proceedings of the 30th international conference on machine learning* (cit. on p. 16).
- Marcus, Marvin (1975). “Finite dimensional multilinear algebra”. In: (cit. on p. 8).
- McCulloch, Warren S and Walter Pitts (1943). “A logical calculus of the ideas immanent in nervous activity”. In: *The bulletin of mathematical biophysics* 5.4, pp. 115–133 (cit. on p. 20).
- Mhaskar, Hrushikesh, Qianli Liao, and Tomaso Poggio (2016). “Learning functions: when is deep better than shallow”. In: *arXiv preprint arXiv:1603.00988* (cit. on p. 18).

- Montufar, Guido F, Razvan Pascanu, Kyunghyun Cho, and Yoshua Bengio (2014). “On the number of linear regions of deep neural networks”. In: *Advances in neural information processing systems*, pp. 2924–2932 (cit. on p. 17).
- Nickolls, John, Ian Buck, Michael Garland, and Kevin Skadron (2008). “Scalable parallel programming with CUDA”. In: *ACM SIGGRAPH 2008 classes*. ACM, p. 16 (cit. on p. 17).
- Orhan, Emin and Xaq Pitkow (2018). “Skip Connections Eliminate Singularities”. In: *International Conference on Learning Representations*. URL: <https://openreview.net/forum?id=HkwBEMWCZ> (cit. on p. 19).
- Pan, Xingyuan and Vivek Srikumar (2016). “Expressiveness of rectifier networks”. In: *International Conference on Machine Learning*, pp. 2427–2435 (cit. on p. 18).
- Pascanu, Razvan, Guido Montufar, and Yoshua Bengio (2013). “On the number of response regions of deep feed forward networks with piece-wise linear activations”. In: *arXiv preprint arXiv:1312.6098* (cit. on p. 17).
- Poggio, Tomaso, Fabio Anselmi, and Lorenzo Rosasco (2015). *I-theory on depth vs width: hierarchical function composition*. Tech. rep. Center for Brains, Minds and Machines (CBMM) (cit. on p. 18).
- Poole, Ben, Subhaneil Lahiri, Maithra Raghu, Jascha Sohl-Dickstein, and Surya Ganguli (2016). “Exponential expressivity in deep neural networks through transient chaos”. In: *Advances in neural information processing systems*, pp. 3360–3368 (cit. on p. 18).
- Raghu, Maithra, Ben Poole, Jon Kleinberg, Surya Ganguli, and Jascha Sohl-Dickstein (2016). “On the expressive power of deep neural networks”. In: *arXiv preprint arXiv:1606.05336* (cit. on p. 18).
- Rumelhart, David E, Geoffrey E Hinton, and Ronald J Williams (1985). *Learning internal representations by error propagation*. Tech. rep. California Univ San Diego La Jolla Inst for Cognitive Science (cit. on pp. 15, 21).

- Simonyan, Karen and Andrew Zisserman (2014). “Very deep convolutional networks for large-scale image recognition”. In: *arXiv preprint arXiv:1409.1556* (cit. on p. 15).
- Srivastava, Nitish, Geoffrey E Hinton, Alex Krizhevsky, Ilya Sutskever, and Ruslan Salakhutdinov (2014). “Dropout: a simple way to prevent neural networks from overfitting.” In: *Journal of Machine Learning Research* 15.1, pp. 1929–1958 (cit. on p. 29).
- Szegedy, Christian, Wei Liu, Yangqing Jia, Pierre Sermanet, Scott Reed, Dragomir Anguelov, Dumitru Erhan, Vincent Vanhoucke, Andrew Rabinovich, et al. (2015). “Going deeper with convolutions”. In: *Conference on Computer Vision and Pattern Recognition* (cit. on pp. 15, 19).
- Van Den Oord, Aaron, Sander Dieleman, Heiga Zen, Karen Simonyan, Oriol Vinyals, Alex Graves, Nal Kalchbrenner, Andrew Senior, and Koray Kavukcuoglu (2016). “Wavenet: A generative model for raw audio”. In: *arXiv preprint arXiv:1609.03499* (cit. on p. 19).
- Weinstein, Alan (1996). “Groupoids: unifying internal and external symmetry”. In: *Notices of the AMS* 43.7, pp. 744–752 (cit. on pp. 45, 84).
- Widrow, Bernard and Marcian E Hoff (1960). *Adaptive switching circuits*. Tech. rep. STANFORD UNIV CA STANFORD ELECTRONICS LABS (cit. on p. 21).
- Wikipedia, contributors (2018a). *Cayley graph* — *Wikipedia, The Free Encyclopedia*. [Online; accessed April-2018]. URL: https://en.wikipedia.org/wiki/Cayley_graph (cit. on p. 53).
- (2018b). *Feedforward neural network* — *Wikipedia, The Free Encyclopedia*. [Online; accessed April-2018]. URL: https://en.wikipedia.org/wiki/Feedforward_neural_network (cit. on p. 15).
- Williamson, S Gill (2015). “Tensor spaces-the basics”. In: *arXiv preprint arXiv:1510.02428* (cit. on p. 8).
- Zell, Andreas (1994). *Simulation neuronaler netze*. Vol. 1. Addison-Wesley Bonn (cit. on p. 15).

Zoph, Barret and Quoc V Le (2016). “Neural architecture search with reinforcement learning”. In: *arXiv preprint arXiv:1611.01578* (cit. on p. 15).

Chapter 0

Trash bin and more drafts

... that we may use in some section.

TODO: Rework 1.1

0.1 Naming conventions

0.1.1 Basic notions

Let's recall the naming conventions of basic notions.

A *function* $f : E \rightarrow F$ maps objects $x \in E$ to objects $y \in F$, as $y = f(x)$.

Its *definition domain* $\mathcal{D}_f = E$ is the set of objects onto which it is defined.

We will often just use the term *domain*.

We also say that f is *taking values* in its *codomain* F .

The *image per* f of the subset $U \subset E$, denoted $f(U)$, is $\{y \in F, \exists x \in U, y = f(x)\}$.

The *image of* f is the image of its domain. We denote \mathcal{I}_f .

A vector space E , which we will always assume to be finite-dimensional in our context, is defined as \mathbb{R}^n , and is equipped with pointwise addition and scalar multiplication.

A *signal* s is a function taking values in a vector space. In other words, a signal can also be seen as a *vector* with an *underlying structure*, where the vector is composed from its image, and the underlying structure is defined by its *domain*.

For example, images are signals defined on a set of pixels. Typically, an image s in RGB representation is a mapping from pixels p to a 3d vector space, as $s_p = (r, g, b)$.

TODO?: figure

0.1.2 Graphs and graph signals

TODO: more defs on grid graphs and other graphs

A *graph* $G = (V, E)$ is defined as a set of nodes V , and a set of edges $E \subseteq \binom{V}{2}$. The words *node* and *vertex* will be used equivalently, but we will rather use the first.

A *graph signal*, or *graph-structured signal* is a signal defined on the nodes of a graph, for which the underlying structure is the graph itself. A *node signal* is a signal defined on a node, in which case it is a *node embedding* in a vector space.

Although this is rarely seen, a signal can also be defined on the edges of a graph, or on an edge. We then coin it respectively *dual graph signal*, or *edge signal* / *edge embedding*.

Graph-structured data can refer to any of these type of signals.

0.1.3 Data and datasets

A dataset of signals is said to be *static* if all its signals share the same underlying structure, it is said to be *non-static* otherwise.

For image datasets, being non-static would mean that the dataset contains images of different sizes or different scales. For graph signal datasets, it would mean that the underlying graph structures of the signals are different.

The point in specifying that objects of a dataset of a machine learning task are signals is that we can hope to leverage their underlying structure.

TODO: figure

0.2 Disambiguation of the subject

This thesis is entitled *Deep learning models for data without a regular structure*. So either the data of interest in this manuscript do not have any structure, or either their structure is not regular.

0.2.1 Irregularly structured data

By structured data, we mean that there exists an underlying structure over which the data is defined. This kind of data are usually modeled as signals defined over a domain. These domains are then composed of objects that are related together by some sort of structural properties. For example, pixels of images can be seen as located on a grid with integer spatial coordinates (a 2d cartesian grid graph).

It then come in handy to define the notions of structure and regularity with the help of graph signals.

Definition 79. Structure

Let $s : D \rightarrow F$ be a signal defined over a finite domain.

An *underlying structure* of the signal s is a graph G that has the domain of s for nodes.

A dataset is said to be *structured*, if its objects can be modeled as signals with an underlying structure.

It is said to be *static* if all its objects share the same underlying structure, and *non-static* otherwise.

In other words, we chose to define “structured data” as “graph-structured data” by some graph. Hence we need to specify for which graphs this structure would be said to be regular, and for which it would not.

Definition 80. Regularity

An underlying structure is said to be *regular*, if it is a regular grid graph. It is said to be *irregular* otherwise.

A dataset is said to be *regularly structured*, if the underlying structures of its objects are regular. It is said to be *irregularly structured* otherwise.

TODO: examples

0.2.2 Unstructured data

Data can also be unstructured. If the data is not yet embedded into a finite dimensional vector space, then we will be interested in embedding techniques used in representation learning. In the other case, it is often possible to fall back to the case of irregularly structured data. For example, vectors can be seen as signals defined over the canonical basis of the vector space, and the vectors of this basis can be related together by their covariances through the dataset. It is typical to use the graph structure that has the canonical basis for nodes, with edges obtained by covariance thresholding.

TODO: examples

What follows is a draft

0.3 Datasets

0.4 Tasks

0.5 Goals

0.6 Invariance

In order to be observed, invariances must be defined relatively to an observation. Let's give a formal definition to support our discussion.

...

0.7 Methods

0.8 Expressivity analysis of dense versus sparse connectivity

Let consider a tensor input x of a neural network layer l . Without loss of generality, we consider that x is a matrix of shape $n \times p$. Its rows are supposed structured by a graph $G = \langle V, E \rangle$, with $|V| = n$, its columns are its feature maps.

In what follows, we discuss the expressivity and efficiency of a dense layer with x as input versus a layer that would leverage G . We start with the regular case and continue onto non-regular structures.

0.8.1 Strong regular case

In the strong regular case, G is a lattice graph such that a convolution is defined naturally on it. For example, this is the case where rows of x defines ticks of a time series, or flattened pixels of an image.

Let consider a convolutional layer $c = (g_c, h_c)$ with padding, defined by q filters of width k . Define y_c its output of shape $n \times q$.

We are interested in knowing if there exists a dense layer that can efficiently replicate c .

Its connectivity matrix W_c is of shape $npxnq$. Obviously, the function g_c can be replicated by a dummy dense layer $d = (g_c, h_d)$ through W_c . However, whereas c has only kpq weights, d has n^2pq . If we consider the families of neural networks \mathcal{C} , \mathcal{D} spanned by their weights θ_c , θ_d , then we realize the \mathcal{C} is less expressive, but in the same time it is more efficient at representing its functions.

Let's define the notion of partial expressivity with respect to a family of functions.

Let \mathcal{F} a family of functions, \mathcal{L} a family of layer functions, and ϵ the approximation coefficient. For $f \in \mathcal{F}$, define $S_\epsilon(\mathcal{L}, f) = \{l \in \mathcal{L}, d(l, f) < \epsilon\}$ and

$$S_\epsilon(\mathcal{L}, \mathcal{F}) = \bigcup_{f \in \mathcal{F}} S_\epsilon(\mathcal{L}, f).$$

TODO: reword above

By abusing and anticipating future correction of this manuscript, we consider that \mathcal{C} and \mathcal{D} are vector spaces. We are interesting in 1. proving that $S_\epsilon(\mathcal{C}, \mathcal{F})$ and $S_\epsilon(\mathcal{D}, \mathcal{F})$ are also vector spaces, and 2. analysing for which \mathcal{F} , $\frac{\dim(S_\epsilon(\mathcal{C}, \mathcal{F}))}{\dim(S_\epsilon(\mathcal{D}, \mathcal{F}))}$ is maximized.

Obviously 1. is false, so 2. is ill-posed (this draft is to be reworded afterward). Instead of using *dim*, we should rather use *card*. However they are potentially infinite families so we should rather use a notion of volume, except if we discretize. So let's discretize.

By the way, "modified" 2. is trivially maximized for $\mathcal{F} = \mathcal{C}$ (and then the ratio equals 1), so let's weaken \mathcal{F} and say it's any family with translation equivariance. We are then interested in proving that if \mathcal{F} is the family on translation equivariant function (on this domains that has to be specified when rewriting this section), then $\frac{\text{card}(S_\epsilon(\mathcal{C}, \mathcal{F}))}{\text{card}(S_\epsilon(\mathcal{D}, \mathcal{F}))}$ is close to 1. Equivariant in our context means commuting with translations (we should rather use the latter expression btw).

The result might be obtained without discretizing as convolutions with padding commutes with translations. Let's guess that they are close to other commutators. In fact that is even it. Proof with Fourier analysis.

0.8.2 Draft

The only dense layer that replicate g_c is obtained through the connectivity matrix W_c . \mathcal{D} is more expressive, however less efficient as we are looking for equivariant functions. It happens that equivariant functions are exactly convolutions with padding.

0.9 Conv drafts

TODO: point

In particular, we have

$$\begin{aligned}
 \forall s \in \mathcal{S}(\Gamma), \tilde{\varphi}(s) &= \tilde{\varphi} \left(\sum_{g \in \Gamma} s[g] \delta_g \right) \\
 &= \sum_{g \in \Gamma} s[g] \tilde{\varphi}(\delta_g) \\
 &= \sum_{g \in \Gamma} s[g] \delta_{\varphi(g)} \\
 &= \sum_{v \in V} s[\varphi^{-1}(v)] \delta_v \\
 \tilde{\varphi}(s) &= \sum_{v \in V} \tilde{\varphi}(s)[v] \delta_v
 \end{aligned}$$

So $\tilde{\varphi}(s)[v] = s[\varphi^{-1}(v)]$ and $\tilde{\varphi}(s)[\varphi(g)] = s[g]$. Let's simplify the notations with $\tilde{\varphi}(s) = t$ and $\varphi(g) = v$, *i.e.* $t[v] = s[g]$ as expected. We then define the group convolution on $\mathcal{S}(V)$ as

$$\begin{aligned}
 (t_1 * t_2)[v] &= (s_1 * s_2)[g] \\
 &= \sum_{h \in \mathcal{G}} s_1[h] s_2[h^{-1}g] \\
 &= \sum_{u \in V} s_1[\varphi^{-1}(u)] h_u(s_2)[\varphi^{-1}(v)] \\
 &= \sum_{u \in V} t_1[u] \tilde{\varphi}(h_u(s_2))[v]
 \end{aligned}$$

$$(t_1 * t_2)[v] = \sum_{u \in V} t_1[u] h_u(t_2)[v] \quad (1)$$

$$(2)$$

TODO: stop sign

Recall that

$$\begin{aligned} \delta_g[h] &= \begin{cases} 1 & \text{if } h = g \Leftrightarrow \varphi(h) = \varphi(g) \\ 0 & \text{otherwise} \end{cases} \\ &= \delta_{\varphi(g)}[\varphi(h)] \end{aligned}$$

$$s = \sum_{v \in V} s[v] \delta_v$$

TODO: lemme on existence of uncountable linearly independent irrational family ?

Proposition 81. The group convolution on $\mathcal{S}(\Gamma)$ has a unique neutral element which is the dirac signal on the identity tranformation.

Proof. Denote δ a neutral element for the group convolution. Note as because of the commutativity the group convolution, a left neutral element is also a right neutral element. We have

$$s[h] = (\delta * s)[h] = \sum_{g \in \Gamma} \delta[g] s[g^{-1}h]$$

which is true for any real valued signal. By chosing a signal π having linearly

independant irrational entries (and using the axiom of choice in case G is not finite), we obtain that

$$\delta[g] = \begin{cases} 1 & \text{if } g = \text{Id} \\ 0 & \text{otherwise} \end{cases} \quad i.e. \quad \delta = \delta_{\text{Id}}$$

Conversely, $(\delta_{\text{Id}} * s)[h] = 1.s[\text{Id}^{-1}h] = s[h]$. \square

In other therms, if there is an isomorphism between Γ and V , the group structure pass to V as well as the definition of the group convolution.

To alleviate this issue, let's introduce the neutral elements δ of the convolution, and the neutral element $\text{Id} \in \Phi^*(V)$.

With the help of δ , we follow the same process as in the proof of Proposition 75, see (2.2), to construct the class of group convolutional operators which defines exactly the class of linear transformations that are equivariant to a certain group.

On graphs, this could be used provided we defined meaningful translations beforehand (see Section ??). Another possibilty would be to search for invariances with respect to graph equivariances and derive a convolution operator similarly than for translations. This approach, which uses group convolutions (Weinstein, 1996), has already been discussed on regular domain to extend CNNs to other invariances than translational ones (Cohen and Welling, 2016; Hoogeboom et al., 2018), as well as on spherical domain with rotation equivariant CNNs (Cohen et al., 2018b). As stated from the previous remark, the big advantage of this approach is that there is no loss of expressivity. However on graphs, this would be more challenging as it's not likely there exists transformations with equivariances. However, let's suppose we found such a set of transformations on a graph, then for Proposition 75 to hold (instead as for regular translations), we see in the proof that they need to be bijective (2.1) and vertex dependent 2.2.

0.9.1

Definition 82. Grounded set of transformations

A set of transformations over a graph $G = \langle V, E \rangle$, *grounded* on a vertex $v_0 \in V$, denoted $\mathcal{P}_{v_0} \subset \Phi(V)$, is a set that is in one-to-one correspondence with V , such that $\forall v \in V, \exists! p_v \in \mathcal{P}_{v_0}, p_v(v_0) = v$.

We have $\mathcal{P}_{v_0} = \text{order}(G) \in \mathbb{N} \cup \{+\infty\}$. For notational convenience we drop the subscript v_0 in what follows.

Definition 83. \mathcal{P} -equivariant convolution operator

Let $G = \langle V, E \rangle$ a graph, not necessarily a grid. Let \mathcal{P} a grounded set of transformations. Then, the \mathcal{P} -equivariant convolution operator f_w is defined as

$$\forall s \in \mathcal{S}(V), f_w(s) = s *_{\mathcal{P}} w = \sum_v s[v] p_v(w)$$

Claim 84. Characterization of \mathcal{P} -eq. convolution operator

The class of linear graph signal transformations that are equivariant to a grounded set \mathcal{P} is exactly the class of \mathcal{P} -equivariant convolutive operations.

Proof. By construction of \mathcal{P} -equivariant convolutions, the proof is similar to the one of Proposition 75. \square

Contents

Temptative previsional plan I	7
1 Introduction	9
1.1 Plan, vision, etc	9
1.2 Deep learning and history	9
1.3 Regular deep learning	9
1.4 Irregular deep learning	9
1.5 Unstructured deep learning	9
1.6 Propagational point of view	9
2 Presentation of the domain	9
2.1 Typology of data	9
2.2 Standardized terminology	9
2.3 Motivation	9
2.4 Datasets	9
2.5 Unifying framework (tensorial product)	9
2.6 Other Unifying frameworks	9
3 Review of models and propositions	9
3.1 How to compare models	9
3.2 Spectral models	9
3.3 Non-spectral	9
3.4 Non-convolutional	9

3.5	Recap and (big) comparison table	9
3.6	Explaining current SOA, current issues, and further work . . .	9
4	Transposing the problem formulation: Structural learning	9
4.1	Structural Representation	9
4.2	Feature visualization (viz on input)	9
4.3	Propagated Signal visualization (viz on S)	9
4.4	Temptatives on learning S	9
4.5	Temptatives on learning S (other)	9
4.6	Covariance-based convolution	9
4.7	Conclusion	9
5	Industrial applications	9
5.1	Context	9
5.2	The Warp 10 platform and Warpscript language	9
5.3	Presentation of use cases: uni vs multi-variate, spatial vs geo, etc	9
5.4	Review and application on regularly structured (spatial) time series	9
5.5	Application to time series database (unstructured)	9
5.6	Application to geo time series (unstructured)	9
5.7	Application to visualization	9
5.8	Market reality (what clients need, what they don't know that can be done ...)	9
5.9	Conclusion	9
6	Conclusion	9
6.1	Summary	9
6.2	Lesson learned	9
6.3	Further avenues	9

<i>CONTENTS</i>	3
Temptative previsional plan II	9
0 Introduction	11
0.1 Teaser	11
0.2 Motivation	11
0.3 Difficulties	11
0.4 Outline	11
1 Presentation of the subject	11
1.1 Title disambiguation	11
1.2 Deep learning	11
1.3 Signals, features, structure, underlying graph	11
1.4 Regular, Irregular, Unstructured	11
1.5 Goals	11
2 Presentation of the field	11
2.0 (<i>Possibly merged with previous chapter</i>)	11
2.1 Tensors	11
2.2 Neural networks	11
2.3 Graphs	11
3 Neural networks on graphs	11
3.1 From non-structured data to graphs	11
3.2 Types of learning	11
3.3 Propagational abstraction	11
3.4 Dense versus sparse connectivity	11
4 Classification of signals over a graph	11
4.0 (<i>Can be splitted in multiple chapters</i>)	11
4.1 Principles	11
4.2 SotA review	11
4.3 Generalizing the convolution	11

4.4	A convolution based on graph translations	11
4.5	Learning the weight sharing scheme	11
4.6	Learning through the covariance matrix	11
4.7	Discussion	11
5	Classification and representation of nodes and graphs	11
5.0	<i>(Optional chapter) either develop this chapter, or just stress out in 3.2 that it won't be</i>	11
5.1	Semi-supervised learning of nodes	11
5.2	Representation learning of graphs	11
5.3	Supervised learning of graphs	11
5.4	Discussion	11
6	Industrial applications	11
6.1	Context, Warp10, (Geo) Time series	11
6.2	Supervised learning	11
6.3	Semi-supervised learning	11
6.4	Representation learning	11
6.5	Market reality and perspectives	11
7	Conclusion	11
7.1	Summary	11
7.2	Discussion	11
7.3	Further avenues	11

<i>CONTENTS</i>	5
-----------------	---

Temptative previsional plan III	11
--	-----------

0 Introduction	13
-----------------------	-----------

0.1 Teaser	13
0.2 Motivation	13
0.3 Difficulties	13
0.4 Outline	13

1 Presentation of the field	13
------------------------------------	-----------

1.1 Tensors	13
1.2 Neural networks	13
1.3 Graphs	13
1.4 Neural networks on graphs	13
1.4.1 Tasks	13
1.4.2 Structures	13
1.4.3 Reusability	13
1.4.4 Approaches	13

2 Classification of signals over a graph	13
---	-----------

2.1 Principles	13
2.1.1 Equivariant functions	13
2.1.2 Partially connected layer	13
2.1.3 About depth	13
2.2 Generalizing the convolution	13
2.3 A convolution based on graph translations	13
2.4 Learning the weight sharing scheme	13
2.5 Other avenues	13
2.6 Discussion	13

3 Classification and representation of nodes and graphs	13
--	-----------

3.0 <i>(Optional chapter) either develop this chapter, or just stress out in 1.4 that it won't be</i>	13
---	----

3.1	Semi-supervised learning of nodes	13
3.2	Representation learning of graphs	13
3.3	Supervised learning of graphs	13
3.4	Discussion	13
4	Industrial applications	13
4.1	Context	13
4.2	Examples	13
4.3	Discussion	13
5	Conclusion	13
5.1	Summary	13
5.2	Discussion	13
5.3	Further avenues	13

<i>CONTENTS</i>	7
Temptative previsional plan IV	13
0 Introduction	15
1 Presentation of the field	15
2 Convolution on graph domains	15
3 Neural networks on graphs	15
4 Industrial applications	15
5 Conclusion	15

Chapter 0

Keywords and tentative titles

Index terms— Deep learning, representation learning, propagation learning, visualization, structured, unstructured regular, irregular, covariant, invariant, equivariant, tensor, scheme, weight sharing, graphs, manifold, euclidean, signal processing, graph signal processing, time series, time series database, distributed application, spatial-time series, geo time series, industrial applications, warp 10, warpscript, ...

Tentative titles

- Learning propagational representations of irregular and unstructured data
- Learning representations of unstructured or irregularly structured datasets
- Propagational learning of unstructured or irregularly structured datasets
- Learning tensorial representation of irregular and unstructured data
- Tensorial representation of propagation in deep learning for irregular and unstructured dataset
- Structural representation learning for irregular or unstructured data

- Word for both “irregularly structured” + “unstructured” = ? (maybe “unorthodox” ?)
- Unorthodox deep learning
- ...
- Deep learning of unstructured or irregularly structured datasets
- Deep learning models for data without a regular structure
- On structures in deep learning
- On deep learning for when data is lacking a regular structure
- Deep learning for non regularly structured data