

Energy-based Multi-Modal Attention

AURELIEN WERENNE



Master Thesis
2018-2019



Energy-based Multi-Modal Attention

Author:
Aurélien WERENNE

Supervisor:
Dr. Raphaël MARÉE

*A thesis submitted in partial fulfillment of the requirements
for the degree of Master of Computer Science and Engineering*

Montefiore Institute
Faculty of Applied Sciences
University of Liège
Liège, Belgium

Academic Year 2018 - 2019

“Sometimes it seems as though each new step towards Artificial Intelligence, rather than producing something which everyone agrees is real intelligence, merely reveals what real intelligence is not.”

Douglas Hofstadter

Abstract

Lorem ipsum dolor sit amet, consectetur adipiscing elit. Sed nec ultrices felis, eu suscipit orci. Nunc laoreet odio velit. Mauris tincidunt quam at purus vulputate vehicula. Etiam eu purus lorem. Sed hendrerit condimentum tellus, ut posuere libero laoreet nec. Sed maximus, tortor aliquam mollis blandit, lectus urna efficitur eros, at laoreet ante turpis suscipit quam. Phasellus eget sollicitudin felis. Sed enim nunc, rutrum vel velit ut, elementum tempor magna. Nullam ut tincidunt orci, et interdum lectus. Proin sed imperdiet tellus. Donec pharetra feugiat leo at fringilla.

Suspendisse consectetur maximus augue. Etiam eu tempor ipsum. Phasellus tempor at purus non cursus. Sed non quam vitae mi rutrum sodales vel posuere odio. Nunc elit arcu, finibus sit amet euismod et, aliquet vel mauris. Mauris eget enim lacus. Donec feugiat eget neque vitae dictum. Nullam turpis neque, mollis at dui quis, lobortis molestie quam. Sed faucibus arcu in odio venenatis, et eleifend lacus lobortis. Pellentesque tincidunt ante non mauris molestie efficitur. Nullam porta massa nulla, at lobortis ex pulvinar sed. Donec auctor consectetur ante, vitae vehicula odio gravida nec. Fusce arcu leo, imperdiet vel magna eu, ullamcorper lacinia orci.

Aliquam vulputate magna lectus, at consequat ante congue et. Sed congue ullamcorper erat, ut volutpat libero consectetur sit amet. Nulla gravida ullamcorper odio, in convallis quam congue sit amet. Nullam tempor ullamcorper pulvinar. Pellentesque rutrum massa eu massa mattis iaculis. Fusce rhoncus tortor id est cursus interdum. Morbi urna elit, commodo sed nisl eget, scelerisque porta dui. Orci varius natoque penatibus et magnis dis parturient montes, nascetur ridiculus mus. Sed nec egestas ligula.

Acknowledgements

Lorem ipsum dolor sit amet, consectetur adipiscing elit. Sed nec ultrices felis, eu suscipit orci. Nunc laoreet odio velit. Mauris tincidunt quam at purus vulputate vehicula. Etiam eu purus lorem. Sed hendrerit condimentum tellus, ut posuere libero laoreet nec.

Thank contributions opens, one day I hope I can do the same...

Contents

Abstract	iii
Acknowledgements	v
1 Introduction	1
1.1 Motivation	1
1.2 Proposed solution	2
1.3 Contributions	3
1.4 Thesis Outline	3
2 Background	5
2.1 Machine Learning	5
2.2 Deep Learning	6
2.3 Physics meets Deep Learning	7
3 Literature Review	9
3.1 Attention in Humans	9
3.2 Attention in Deep Learning	10
4 Energy Estimation	13
4.1 Autoencoder	13
4.2 Energy in Autoencoders	15
4.3 Experiment I	16
4.4 Limitations	17
5 Energy-based Multi-Modal Attention	19
5.1 General Framework	19
5.2 From Potential to Modal energies (step 2)	20
5.3 From Modal energies to Importance scores (step 3)	22
5.4 From Importance to Attention scores (step 4)	22
5.5 Training & Regularization	24
5.6 Advantages	26
6 Experiments & Results	29
6.1 Pulsar detection	29
6.2 Corruption	29
6.3 Experiment II	29
6.4 Experiment III	30
7 A Unified Model for Multi-Modal Attention	33
8 Conclusion	35
8.1 Contributions	35
8.2 Research questions	35

8.3 Future work	35
A Dataset	37
B Miscellaneous	41
B.1 Integrability criterion	41
B.2 Gradient with respect to gamma	41
Bibliography	43

List of Figures

1.1	Lidar & Camera view in self-driving cars	2
1.2	Multi-Modal model with/without EMMA	3
2.1	Venn diagram of the Artificial Intelligence field	5
2.2	Early and late fusion	7
2.3	Energy surface evolution	8
3.1	Looking to Listen framework	11
3.2	Noise-tolerant fusion model	11
4.1	The architecture of the two families of autoencoders	13
4.2	Vectorial representation of undercomplete AE	14
4.3	Vectorial representation of overcomplete AE	14
4.4	Vector field circle manifold	15
4.5	Manifold generation of 200 samples	17
4.6	Vector fields on wave and circle manifold	17
4.7	Heatmap of estimators on wave and circle manifold	18
5.1	High-level view of a Multi-Modal Network with EMMA	19
5.2	Summary of main steps in EMMA	20
5.3	Input-output of Boltzmann distribution for two different temperatures	22
5.4	Attention function	23
5.5	Summary of end-to-end training	27
7.1	A possible architecture for a unified multi-modal attention	33
A.1	Evolutionary endpoints for main sequence stars	38
A.2	Lighthouse model of a radio pulsar	39
A.3	Pulse profiles of two separate pulsars	39
A.4	Signal Dispersion	40

Notation and Acronyms

\triangleq	Is defined as
N	Number of samples
M	Number of modes
k_B	Boltzmann constant
M_\odot	Solar mass
e	Euler's number, base of the natural logarithm (2.71828)
\mathcal{L}	Loss function
$\boldsymbol{\theta}$	Set of parameters of the specified model
$\nabla_{\boldsymbol{\theta}}$	Gradient with respect to $\boldsymbol{\theta}$
λ_c	Weight of capacity penalty
λ_e	Weight of energy penalty
Ω	Energy regularizer
Ψ_i	Potential energy of mode i
E_{total}	Total energy
E_i	Modal energy of mode i
e_i	Self-energy of mode i
e_{ij}	Shared energy of mode j on mode i
α_i	Importance score of mode i
β_i	Attention score of mode i
ρ	Coldness in Boltzmann distribution
T	Temperature in Boltzmann distribution
AE	A utoeconder
BP	B ack-propagation
CNN	C onvolutional N eural N etwork
DAE	D enoising A utoeconder
DL	D eep L earning
DM	D ispersion M easure
EMMA	E nergy-based M ulti- M odal A ttention
ISM	I nterstellar M edium
IP	I ntegrated P rofile
LSTM	L ong S hort T erm M emory
MMDL	M ulti M odal D eep L earning
MMN	M ulti M odal N etwork
NLP	N atural L anguage P rocessing
RNN	R ecurrent N eural N etwork
SGD	S tochastic G radient D escent
SNR	S ignal-to- n oise R atio
WER	W ord E rror R ate

Chapter 1

Introduction

1.1 Motivation

In recent years, there has been tremendous advances in the field of Artificial Intelligence (AI), especially in Deep Learning (LeCun, Bengio, and Hinton, 2015; Fan, Ma, and Zhong, 2019). Deep Learning has helped AI systems reach and sometimes surpass human-level perception, mostly in computer vision (He et al., 2016) and natural language processing (Wu et al., 2016). Giving rise to amazing industrial application such as autonomous driving, early cancer detection, enhanced machine translation, etc. A primary concern of engineers is to make sure the trained models strive to be error-free, which can be challenging if the input data does not carry enough information.

One possible solution researchers started to explore is to use multiple modalities, which makes sense since our experience of the world is multi-modal, i.e., we see objects, hear sound, feel the texture, smell odours, and taste flavours. The term modality, also called mode, is generally understood to mean "the way in which something happened or is experienced" (Baltrušaitis, Ahuja, and Morency, 2019). Multi-Modal Deep Learning (MMDL) is used in the hope that the information carried by each mode is additive, such that the model can learn to make more accurate predictions. For example, in (Caltagirone et al., 2018) sensorial inputs from wide angle cameras and LIDAR¹ sensors are combined for road detection. Because cameras provide dense information over a long range under good illumination and fair weather, whereas LIDARs are only marginally affected by the external lighting conditions and provide accurate distance measurements but have a limited range. Thus, fusing the sensorial information with deep learning gets the advantages of both sensors. Despite its efficacy, MMDL suffers from a major drawback: no explicit mechanisms exist to handle failing modes. In the present report, a mode is said to be failing if a) it contains a significant amount of noise, b) the data is much different from the training data, c) the data is missing. Failing modes a) and b) generally degrade the quality of the predictions because they introduce perturbations in the network.

On the other hand, humans seem to handle these situations robustly on a daily basis. A famous example showing this ability is called the cocktail-party effect (Cocktail party effect, 2010), it refers to the difficulty we sometimes have understanding speech in noisy social settings. As a subconscious response, we tend to look at the mouth of our interlocutor i.e. we shift some attention from the auditory to the visual senses. Similarly, our attention is shifted from vision to touch when we are wandering in a room where the lights suddenly switch off. These examples indicate that humans handle modes with

¹Laser Detection and Ranging

perturbations (first example) or missing information (second example) by shifting their attention on the other more relevant modes (Driver and Spence, 1998).

Inspired from this behaviour, this report presents a new approach to tackle failing modes. I introduce a novel attention mechanism, named *Energy-based Multi-Modal Attention* (EMMA), able to decide how much attention to devote to each mode, such that the relevant information is kept while masking out the perturbations. Additionally, this work offers some important insight into how current attention mechanisms in deep learning are surprisingly similar in some ways to attention in humans.

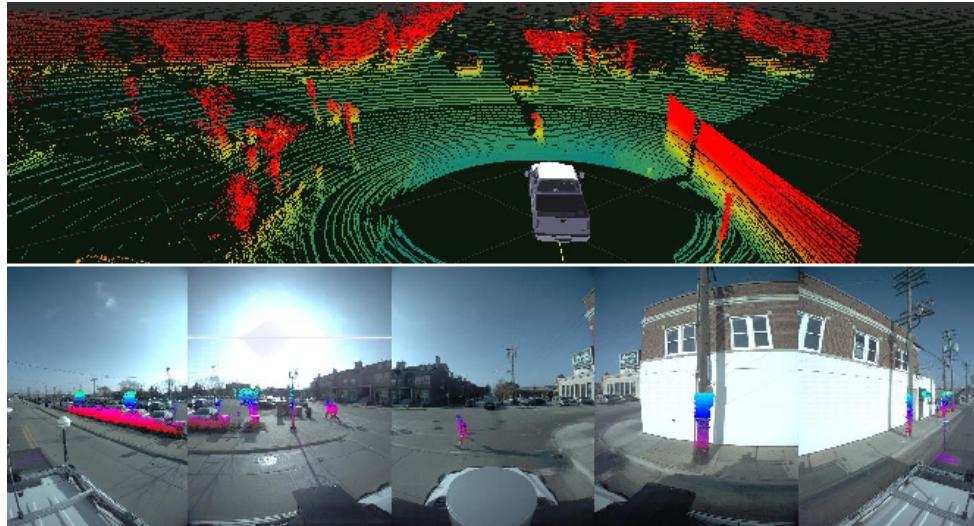


FIGURE 1.1: Same environment, different modes (top: LIDAR view, bottom: camera view)

1.2 Proposed solution

The attention module EMMA is inserted in front of the model, focusing its attention on the modes such that the most useful information passes through while the perturbations are filtered out. The amount of attention distributed to a mode is based on its importance, encompassing three intrinsically tied properties:

- *relevance*: the quantitative influence of the mode over the predictions
- *failure intensity*: a measure proportional to the outlyingness of the mode
- *coupling*: how does the information carried by the mode relate to the other modes? Is it redundant, complementary or conflicting?

Let us emphasize that the determining the importance is sample dependent, and is thus not easily solved by learning the global tendency.

Software Implementation

All the implemented models and experiments are available at this repository², with a wiki explaining how to run the experiments; PyTorch³ was the main framework used

²<https://github.com/Werenne/energy-based-multimodal-attention>

³<https://pytorch.org/>

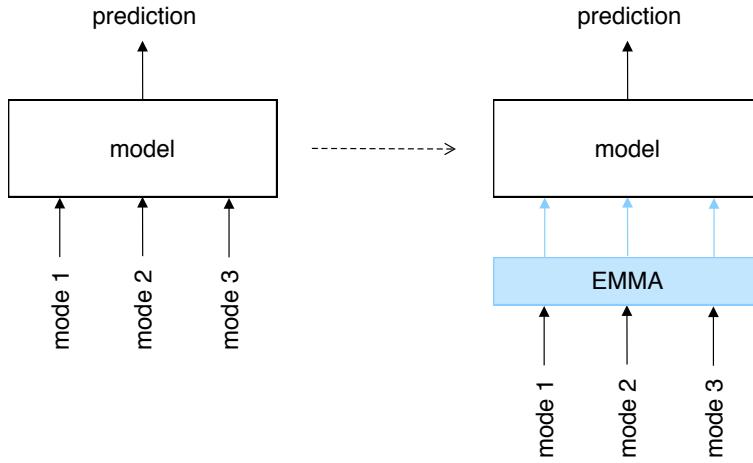


FIGURE 1.2: A multi-modal model with three input modes, without EMMA (left), improved with EMMA (right)

regarding the Machine Learning part.

1.3 Contributions

The work presented in this Master thesis has led to three novel contributions.

Contribution 1: an attention module improving significantly the robustness against failing modes. In Chapter 5, we discuss the design of a new attention mechanism based on energy models (LeCun et al., 2006), that can be added to any multi-modal model. TODO: present results obtained.

Contribution 2: a simple yet powerful regularizer on attention mechanisms.

We slightly modify a common attention function permitting us to establish a link to the concept of capacity in psychology (Kahneman, 1975); Capacity is the amount of attention distributed among the inputs. Subsequently, a new regularizer is introduced to control the capacity, which we claim can help generalize against unexpected situations.

Contribution 3: a unified model for multi-modal attention. In Chapter 3, a review of the literature on attention in humans helps us identify how to construct a more complete multi-modal attention module.

1.4 Thesis Outline

The remainder of this work is organised as follows.

Chapter 2 explains the background (i.e. deep learning and energy models) this work is based upon.

Chapter 3 reviews the literature about attention in psychology and deep learning, and the similarities between them.

Chapter 4 describes a method for the estimation of the failure intensity of a mode.

Chapter 5 presents the ideas and architecture of the Energy-based Multi-Modal Attention module (Contribution 1 & 2).

Chapter 6 presents a thorough evaluation and analysis of the module outlined in Chapter 5.

Chapter 7 proposes a unified multi-modal attention module (Contribution 3).

Chapter 8 concludes this work and suggests possible directions for future research.

Chapter 2

Background

2.1 Machine Learning

Machine Learning is a subfield of Artificial Intelligence (see Figure 2.1) concerned with the design of algorithms that allow machines (e.g. computers, robots, embedded systems) to learn. For a task \mathbf{T} , a performance measure \mathbf{P} and an amount of data \mathbf{D} , the system is said to be learning if it improves its performance \mathbf{P} at the task \mathbf{T} by increasing \mathbf{D} (gain experience). Moreover, there are three main types of learning paradigms, namely supervised, unsupervised and reinforcement learning. In supervised learning (Loog, 2017), the model learns on a labeled dataset, providing an answer that the algorithm can use to evaluate its accuracy on training data. An unsupervised model (Ghahramani, 2004), on the contrary, extracts features and patterns from unlabelled data. Lastly, reinforcement learning (Li, 2017) is typically used to train agents in dynamic environments, where the agent is able to act upon the environment. Reinforcement learning is best explained by an analogy, the learning algorithm is like a dog trainer, which teaches the dog (agent) how to respond to specific signs, like a whistle for example. Whenever the dog responds correctly, the trainer gives a reward to the dog, reinforcing the correct behaviour of the dog. Based on these three paradigms, several families of algorithms were invented. Deep learning (Fan, Ma, and Zhong, 2019) is one of those families and is particularly powerful on perception tasks.

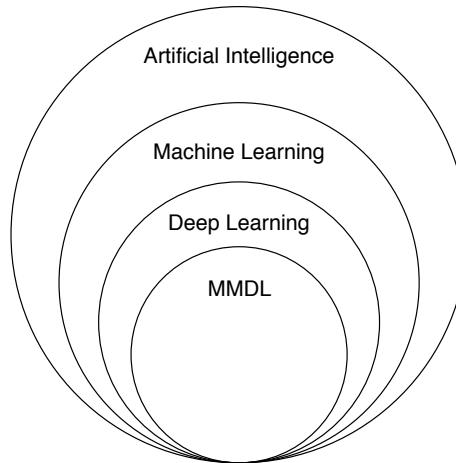


FIGURE 2.1: Venn diagram of the Artificial Intelligence field

2.2 Deep Learning

Deep learning models, also called Deep Neural Networks, offer the significant advantage of being able to learn their own feature representation for the completion of a given task. Neural Networks were loosely inspired from our own brains, but nowadays can best be seen as gigantic composed non-linear parametric functions. The parameters are tuned by optimizing a loss function with Stochastic Gradient Descent (SGD) or one of its many enhancements (Ruder, 2016). Let $\boldsymbol{\theta}$ be the set of parameters, \mathcal{L} the loss function, y the groundtruth (labels) and \hat{y} the predictions. First, the SGD algorithm estimates the gradient of the cost function on a randomly sampled batch of size N as

$$\mathbf{g} = \frac{1}{N} \nabla_{\boldsymbol{\theta}} \sum_{i=1}^N \mathcal{L}(\hat{y}^{(i)}, y^{(i)}) \quad (2.1)$$

where the computation of the gradient itself is done using back-propagation (BP) (Chauvin and Rumelhart, 1995). The SGD algorithm then follows the estimated gradient downhill, $\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} - \epsilon \mathbf{g}$ where ϵ is the learning rate, in the hope of minimizing the loss.

Optimizing the parameters to represent all valid inputs of a task, where the data is often very high-dimensional (e.g., images, sounds, text), may seem hopeless. However, neural networks surmount this obstacle by assuming that these high-dimensional data are lying along low-dimensional manifolds¹ (Goodfellow, Bengio, and Courville, 2016). An intuitive observation in favor of this claim is that uniform noise essentially never resembles structured inputs from these tasks. More rigorous experiments supporting the manifold hypothesis are (Cayton, 2005; Narayanan and Mitter, 2010; Weinberger and Saul, 2006).

Multi-Modal Deep Learning

Remember, a modality refers to "the way in which something happened or is experienced" (Baltrušaitis, Ahuja, and Morency, 2019). Multi-Modal Deep Learning (MMDL) is simply the research area of neural networks using input samples consisting of multiple modes. Baltrušaitis et al. identified five non-exclusive use-cases of MMDL,

- *Representation*: learning how to represent and summarize multimodal data in a way that exploits the complementarity and redundancy
- *Translation*: learning how to map data from one modality to another (e.g., image captioning)
- *Alignment*: learning to identify the direct relationships between elements from two or more different modalities (e.g. alignment of sound and video)
- *Fusion*: learning to join information from two or more modalities to perform predictions
- *Co-learning*: learning to transfer knowledge between modalities and their respective predictive models (e.g., zero shot learning)

The EMMA module is applied to multi-modal networks performing fusion. Furthermore, networks doing fusion can combine their modalities in three different ways: by

¹A manifold designates a connected set of points that can be approximated well by considering only a small numbers of degrees of freedom

early-fusion, late-fusion and an hybrid of the first two. Early-fusion architectures have uni-modal encoders extracting the features of each mode, the obtained features are then concatenated altogether and fed into a common decoder making the predictions (see Figure 2.2a). In contrast, late-fusion has uni-modal predictors for each mode, followed by a decoder weighting the uni-modal predictions to compute the final prediction.

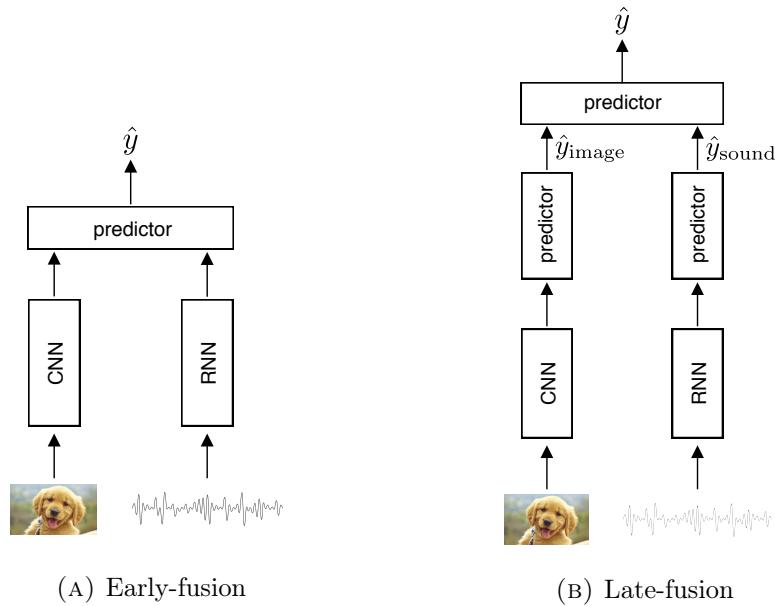


FIGURE 2.2: Fusion of images and sounds for a classification task with a Convolutional Neural Network (CNN) (He et al., 2016), Recurrent Neural Network (RNN) (Wu et al., 2016)

2.3 Physics meets Deep Learning

Modelling complex probability distributions by parametric functions such as deep learning models is a difficult task, because the probabilities must be positive and sum up to one. Back in time, many researchers in deep learning had an academic background in physics, from which they naturally found inspiration. The distribution of kinetic energies among molecules of gaz, called Boltzmann distribution, is given by

$$p(E_i) = \frac{1}{Z} e^{-E_i/k_B T} \quad \text{with the partition function} \quad Z = \int e^{-E_j/k_B T} \quad (2.2)$$

where E_i is the kinetic energy, k_B the Boltzmann constant and T the temperature of the environment. TODO: interpretation. Deep learning researchers used the Boltzmann distribution to construct energy-based models (LeCun et al., 2006), which are more convenient to learn probability distributions. Neural networks of this type learn to model an energy function which can then mapped to the target distribution using Equation (2.2). The learned energy functions are low on the data manifold and high everywhere else (see Figure 2.3), since $E_i \propto -\log p_i$. A few examples of efficient energy-based models are Generative Adversarial Networks (GAN) (Goodfellow et al., 2014), Variational Autoencoders (VAE) (Kingma and Welling, 2013) and Denoising Autoencoders (DAE) (Vincent et al., 2008). A simple energy-based model will be used in this work as a measure of the outlyingness of the data (more in Chapter 4).

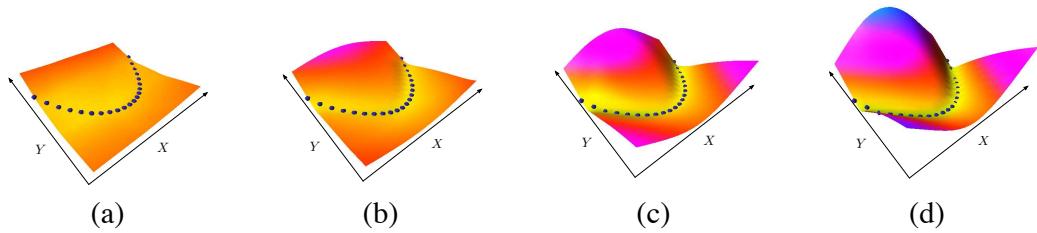


FIGURE 2.3: The shape of the energy surface at four intervals. Along the x-axis is the variable X and along the y-axis is the variable Y . The shape of the surface at (a) the start of the training, (b) after 15 epochs over the training set, (c) after 25 epochs, and (d) after 34 epochs. The energy surface has attained the desired shape: the energies around the training samples are low and energies at all other points are high. *Image and caption from (LeCun et al., 2006)*

Chapter 3

Literature Review

The purpose of this chapter is to review the state-of-the-art literature of multi-modal attention. The first section describes attention in humans from both a psychological and a neurological point of view. We argue this will give the reader more intuition about attention in deep learning. The second part moves on to the different attention mechanisms in deep learning, in particular self-attention and crossmodal attention.

3.1 Attention in Humans

The most profound effect of attention is its capacity to bring the attended stimuli into the forefront of our conscious experience while unattended stimuli fades into the background, increasing the processing efficiency at every stage of perception (Watzl, 2017). A widely held assumption in the psychology literature is that the most fundamental function of attention is selection. But also at the level of single neurons, neuroscientist typically thought of attention in terms of selection between stimuli competing for the same neural receptive field (Desimone and Duncan, 1995). Daniel Kahneman, an authority in psychology and economy, investigated the way in which humans perform multi-tasking (i.e., solve a multi-modal problem). Kahneman claimed that attention was more than selection, that it could be viewed as a limited resource being shared among the different modes, but he could not generalize his findings to the intra-modal level¹. Moreover, the selection theory has been vigorously challenged in recent years by the amplification theory, where attention is an additional activity that interacts with built-in perceptual mechanisms by amplifying some of the input signals (Fazekas and Nanay, 2018). Furthermore, the absolute intensity of amplification is not important, in contrary it is the relative intensity between the inputs that matters (*the contrast effect*). Notice that the amplification theory generalizes the concept of capacity to the intra-modal level and neural level. Interestingly, we will see that the basic principles of attention mechanisms in deep learning has significant similarities with amplification.

Regarding multi-modal attention, three types can be distinguished: endogenous, exogenous and crossmodal attention (Driver and Spence, 1998). People orient their attention endogenously whenever they voluntarily choose to attend to something, such as when listening to a particular individual at a noisy cocktail party, or when concentrating on the texture of the object that they happen to be holding in their hands. By contrast, exogenous orienting occurs when a person's attention is captured reflexively by the sudden onset of an unexpected event, such as when a mosquito suddenly lands

¹Intra-modal attention manifests itself only in a subset of the mode, whereas inter-modal attention is between modes.

on our arm. Lastly, crossmodal attention refers to the interaction of attention between two or more modes such as using visual clues (e.g. lip movements) to focus on the voice of a particular individual at a noisy cocktail party.

3.2 Attention in Deep Learning

Attention-based approaches in deep learning focus on network architectures that specifically attend to regions of their input space. The most common way to do this, is by multiplying the input by an attention mask, where the attention mask consist of normalized continuous values between zero and one. Observe the similarity with the amplification theory described in the previous section. In self-attention (Bahdanau, Cho, and Bengio, 2014), the attention mask is computed from the same mode on which it is applied. Conversely, for crossmodal attention mechanisms (Li et al., 2019), the attention mask is computed from multiple modes.

Self-attention was first introduced in natural language processing (NLP) for machine translation tasks by (Bahdanau, Cho, and Bengio, 2014). It helped the translation task by enabling the model to automatically search for parts of a source sentence that are relevant to predicting the next target word. With this approach, Bahdanau et al. achieved a translation performance comparable to the existing state-of-the-art phrase-based system on the task of English-French translation. Since then it has become a prominent tool in NLP but has also been used in a variety of other tasks such as image classification. (Hoogi et al., 2019) uses self-attention to learn to suppress irrelevant regions in images and highlight salient features useful for the specific classification task. The authors in (Hoogi et al., 2019) reduced the computation load and were able to compensate the absence of a deeper network by using the self-attention, without having a decreased classification performance. For a detailed review on this self-attention mechanisms, see (Galassi, Lippi, and Torroni, 2019).

In (Ephrat et al., 2018), an audio-visual model is presented for isolating a single speech signal from a mixture of sounds such as other speakers and background noise (see Figure 3.1). Crossmodal attention is used to focus on certain parts of the audio with respect to an image of the desired speaker. The authors showed superior results compared to state-of-the-art audio-only methods. Similar works (Libovický, Helcl, and Mareček, 2018; Li et al., 2019; Wang, Wang, and Wang, 2018) are using crossmodal attention and have attained impressive results. However, most research using crossmodal attention has tended to focus on obtaining better predictions rather than improving the robustness. A few exceptions are discussed below.

A work investigating how multimodal fusion can help against failing modes is (Afouras et al., 2018). Their model fuses audio and video to obtain better speech-to-text. Interestingly, Afouras et al. use a combination of self-attention mechanisms followed by a crossmodal attention layer. The model was tested on thousands of natural sentences of British television. Furthermore, they added babble noise with 0dB signal-to-noise ratio to the audio streams, where the babble noise samples are synthesized by mixing the signals of 20 different audio samples from the dataset. The audio-visual model achieved a 13.7% word error rate (WER) on the dataset without noise, and a 33.5% WER on the dataset with noise whereas the audio-only model only achieved 64.7% WER. Despite obtaining great results, a major weakness with this experiment, however, is that their test set is corrupted in the exact same manner as their training set; there are also no varying levels of noise between the two modes. Additionally, the attention module is not

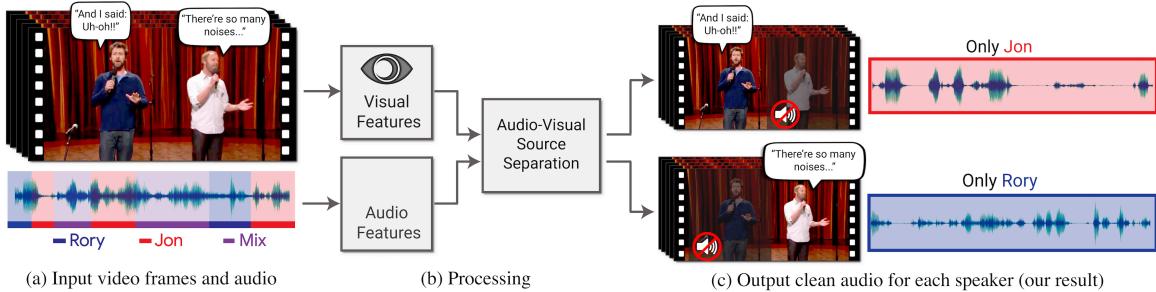


FIGURE 3.1: The authors of (Ephrat et al., 2018) present a model for isolating and enhancing the speech of desired speakers in a video. (a) The input is a video (frames + audio track) with one or more people speaking, where the speech of interest is interfered by other speakers and/or background noise. (b) Both audio and visual features are extracted and fed into a joint audio-visual speech separation model. The output is a decomposition of the input audio track into clean speech tracks, one for each person detected in the video (c). This allows them to then compose videos where speech of specific people is enhanced while all other sound is suppressed. Their model was trained using thousands of hours of video segments from our new dataset, AVSpeech. The ‘Stand-Up’ video (a) is courtesy of Team Coco. *Image and caption from (Ephrat et al., 2018)*

constructed to detect unseen samples. To sum it up, the model was not tested against realistic failing modes situations.

The work that is most relevant to our proposed method is the attentive context proposed in (Shon, Oh, and Glass, 2018), which also incorporates attention on the inter-modal level to filter the perturbations (see Figure 3.2). The model is evaluated on a face-verification task, receiving a voice and a face. The attention mask $[\alpha_v, \alpha_f]$ is computed via a linear function, $f_{att} = \mathbf{W}^T [\mathbf{e}_v, \mathbf{e}_f] + \mathbf{b}$, on the embeddings \mathbf{e}_v and \mathbf{e}_f . Several defects of the attention function f_{att} can be observed:

1. The function is likely to be too simple to neither capture complicated dependencies between the modes, nor to recognize unseen data.
2. The way in which the function is designed forces the model to extract embeddings of the same size, which may be a significant constraint when combining modes from low and high-dimensional data.
3. TODO: explain no control (capacity, temperature)

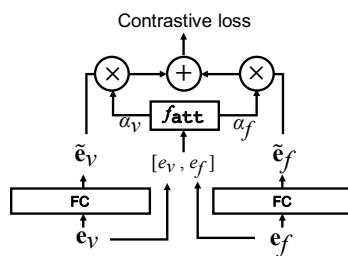


FIGURE 3.2: Neural network based fusion approaches. \mathbf{e}_v : speaker embedding, \mathbf{e}_f : face embedding. FC denotes a fully connected layer.

Image from (Shon, Oh, and Glass, 2018)

Chapter 4

Energy Estimation

As mentioned in the introduction, the EMMA module needs a measure of the outlyingness of each mode, in order to determine which modes are important and which are not. To this end we use an energy-based model (see Section 2), since their energy function is proportional to the negative log-likelihood (NLL). This chapter discusses how an energy function can be derived from an autoencoder.

4.1 Autoencoder

Autoencoders (AE) are models trained to reproduce their inputs to their outputs. An autoencoder is composed of two main parts, the encoder f and the decoder g . The input $\mathbf{x} \in \mathbb{R}^L$ is passed through the encoder as $f(\mathbf{x}) = h(W_f \mathbf{x} + \mathbf{b}_f) = \mathbf{u}$ where $h(\cdot)$ is an activation function and \mathbf{u} represents the hidden layer. The decoder is then in charge of reconstructing the input, $g(\mathbf{u}) = W_g \mathbf{u} + \mathbf{b}_g$. The output is often called the reconstruction and is written $r(\mathbf{x}) = g(f(\mathbf{x}))$. Autoencoders are trained in an unsupervised manner, most of the time using the mean-squared error between input and output as a loss function, $\mathcal{L}_{\text{MSE}} = \|r(\mathbf{x}) - \mathbf{x}\|_2^2$. Training a model to copy its input may seem useless. To answer this point, we need to distinguish two families of autoencoders: the undercomplete and overcomplete autoencoders.

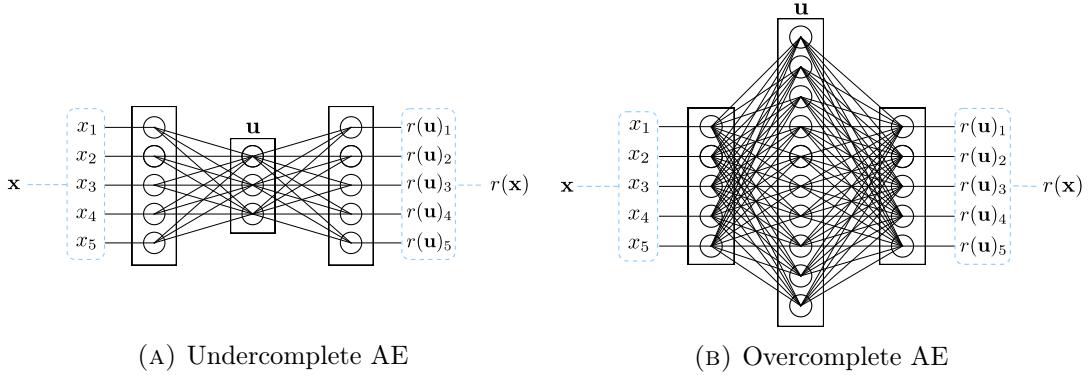


FIGURE 4.1: The architecture of the two families of autoencoders

Undercomplete

An autoencoder is said to be undercomplete, when the size of the hidden layer \mathbf{u} is smaller than the size of the input/output layers (see Figure 4.1a). As a result, the input

\mathbf{x} has to pass through a bottleneck, forcing the model to loose some information and keeping only the most relevant features. It can be thought of as non-linear principal component analysis (Scholz, Fraunholz, and Selbig, 2008; Ladjal, Newson, and Pham, 2019): the values formed in the hidden layer are a non-linear representation in latent space of the input. As can be seen in Figure 4.2, minimizing the mean squared error is similar to minimizing the norm of the vector $r(\mathbf{x}) - \mathbf{x}$.

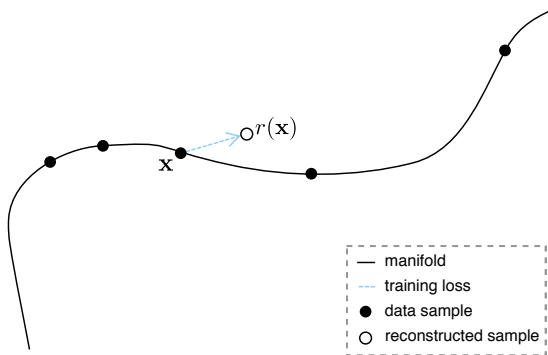


FIGURE 4.2: Vectorial representation of an undercomplete reconstruction process

Overcomplete

Conversely, an overcomplete AE has more hidden units than its input/output layer (see Figure 4.1b). Straightforwardly, the model can thus learn to perfectly copy its input to the output through the L hidden units. To spice things up, the input is corrupted before being passed through the encoder. If we force the AE to reconstruct the original input, we now have a model learning to denoise signals. This type of AE is called a denoising autoencoder (DAE). More formally, the input is corrupted with some small isotropic noise $\tilde{\mathbf{x}} = \mathbf{x} + \epsilon$ where $\epsilon \sim \mathcal{N}(0, \sigma^2)$, with the training loss

$$\mathcal{L}_{\text{MSE}} = \|r(\tilde{\mathbf{x}}) - \mathbf{x}\|_2^2 \quad (4.1)$$

Notice the difference with the loss function of the undercomplete AE. We verify on Figure 4.3 that minimizing the loss, $r(\tilde{\mathbf{x}}) \rightarrow \mathbf{x}$, is equivalent to learning to invert the corruption, $r(\tilde{\mathbf{x}}) - \tilde{\mathbf{x}} \rightarrow -(\tilde{\mathbf{x}} - \mathbf{x})$.

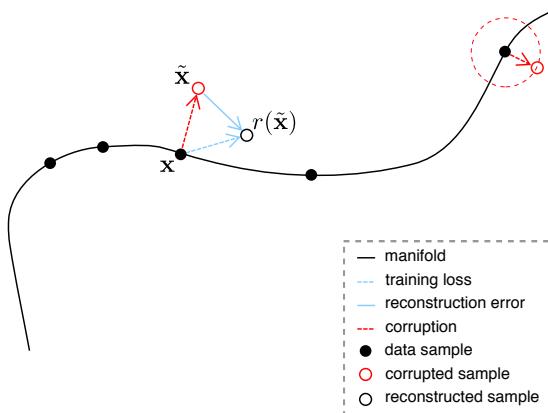


FIGURE 4.3: Vectorial representation of an overcomplete reconstruction process

4.2 Energy in Autoencoders

The authors in (Alain and Bengio, 2012) found that the reconstruction error of a trained denoising autoencoder is proportional to the score (gradient of log-likelihood)

$$r(\tilde{\mathbf{x}}) - \tilde{\mathbf{x}} \propto \frac{\partial \log p(\tilde{\mathbf{x}})}{\partial \tilde{\mathbf{x}}} \quad (4.2)$$

To put it differently, the reconstruction error points towards the corresponding most likely datapoint. This result is not particularly surprising, indeed, denoising a signal is essentially equivalent to finding the most likely datapoint in the nearby neighborhood (see Figure 4.3). To illustrate Equation (4.2), we train a DAE on a generated circle manifold (more details about this experiment in Section 4.3). As we can see below, the vector field of the reconstruction error does indeed point towards the data manifold.

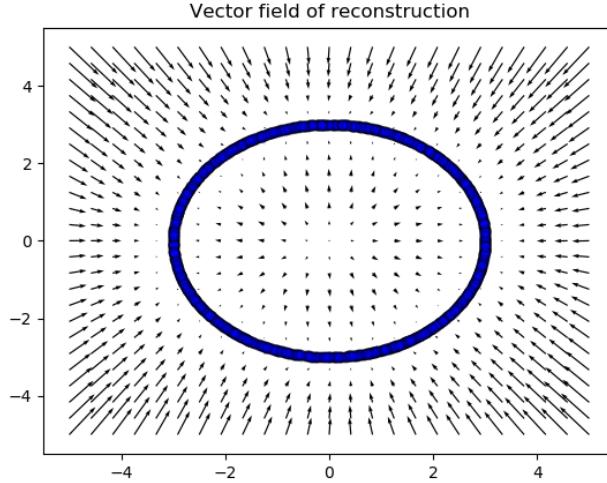


FIGURE 4.4: Vector field of reconstruction error on circle manifold.
No corruption is applied at test time, the reconstruction error vector is simply the output minus the input.

In (Kamyshanska and Memisevic, 2014), authors observed that using tied weights ($W_f = W_g^T$), turns the integrability criterion¹ satisfied:

$$\begin{aligned} \frac{\partial(r(\tilde{\mathbf{x}})_i - x_i)}{\partial x_j} &= \sum_k W_{ik} \frac{\partial h(W\tilde{\mathbf{x}} + \mathbf{b}_f)}{\partial(W\tilde{\mathbf{x}} + \mathbf{b}_f)} W_{jk} - \delta_{ij} \\ &= \frac{\partial(r(\tilde{\mathbf{x}})_j - x_j)}{\partial x_i} \end{aligned} \quad (4.3)$$

where δ_{ij} denotes the Kronecker delta, $W = W_f$ and W_{ij} is the element on the i^{th} row and j^{th} columns. The vector field under those circumstances can be expressed as a gradient of a scalar field $-\Psi$, such that $r(\tilde{\mathbf{x}}) - \tilde{\mathbf{x}} = -\nabla\Psi(\tilde{\mathbf{x}})/\partial\tilde{\mathbf{x}}$. In analogy to physics, the vector field can be interpreted as a force applied on the input and the scalar field as a potential energy. Thereupon, the reconstruction process can be seen as a gradient descent in the potential energy landscape (Kamyshanska and Memisevic, 2014). For our purpose, an important observation to make is that the potential energy is proportional

¹See Appendix B.1

to the NLL,

$$\frac{\partial \Psi(\tilde{\mathbf{x}})}{\partial \tilde{\mathbf{x}}} \propto -\frac{\partial \log p(\tilde{\mathbf{x}})}{\partial \tilde{\mathbf{x}}} \Rightarrow \Psi \propto -\log p \quad (4.4)$$

The potential energy being the gradient of the reconstruction error, we can compute Ψ as

$$\Psi(\tilde{\mathbf{x}}) = - \int (r(\tilde{\mathbf{x}}) - \tilde{\mathbf{x}}) d\tilde{\mathbf{x}} \quad (4.5)$$

Substituting $f = h(W\mathbf{x} + \mathbf{b}_f)$ and $g = W^T\mathbf{x} + \mathbf{b}_g$

$$\Psi(\tilde{\mathbf{x}}) = - \int f(\tilde{\mathbf{x}}) d\tilde{\mathbf{x}} + \frac{1}{2} \|\tilde{\mathbf{x}} + \mathbf{b}_r\|_2^2 + \text{const} \quad (4.6)$$

In this work only the sigmoid activation function will be used, thus solving f

$$\Psi(\tilde{\mathbf{x}}) = - \sum_k \log(1 + \exp(W_k^T \tilde{\mathbf{x}} + b_k^f)) + \frac{1}{2} \|\tilde{\mathbf{x}} + \mathbf{b}_r\|_2^2 + \text{const} \propto -\log p(\tilde{\mathbf{x}}) \quad (4.7)$$

where W_k^T is the k^{th} column of W^T , and b_k^f the k^{th} element of \mathbf{b}_f . The intermediate steps between (4.5) and (4.7) are detailed in (Kamyshanska and Memisevic, 2014). In contrast to physics, notice that the potential energy can be negative by construction.

4.3 Experiment I

In this experiment, two simple data manifolds are generated, on which separate denoising autoencoders are trained. From those trained autoencoders, the energy function is computed on a grid mesh. As a comparison, we also compute the reconstruction error, $\|r(\tilde{\mathbf{x}}) - \tilde{\mathbf{x}}\|_2^2$, which is sometimes used in the Machine Learning community as a way to detect outliers.

Manifolds

The manifolds consists of a set of N samples $\mathbf{x} \in \mathbb{R}^2$ in the form of a wave and a circle. The N samples, written \mathbf{t} , are randomly selected in an interval $[0, 2\pi]$, and are transformed to manifolds as

$$\begin{array}{ll} \text{wave} & \begin{cases} \mathbf{x}_1 = \mathbf{t} - \pi \\ \mathbf{x}_2 = \sin(\mathbf{t}) \end{cases} \\ & \text{circle} \begin{cases} \mathbf{x}_1 = 3 \sin(\mathbf{t}) \\ \mathbf{x}_2 = 3 \cos(\mathbf{t}) \end{cases} \end{array}$$

The result of this process can be viewed on Figure 4.5.

Setup

Each autoencoder has 8 hidden units, is trained for 25 epochs, with a batch size of 100, a corruption noise $\sigma = 0.008$ and a learning rate of $1e^{-3}$. The used optimizer is *Adam* (Kingma and Ba, 2014).

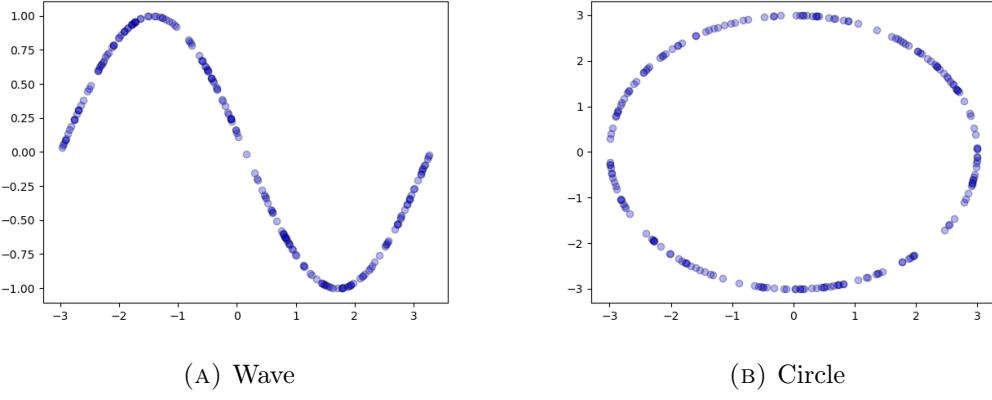


FIGURE 4.5: Manifold generation of 200 samples

Results

As expected the vector fields of the reconstruction error are directed towards the manifolds (see Figure 4.6), the manifolds acting as sinks in the vector field. Notably, observe the presence at the origin for the circle manifold (see Figure 4.6b).

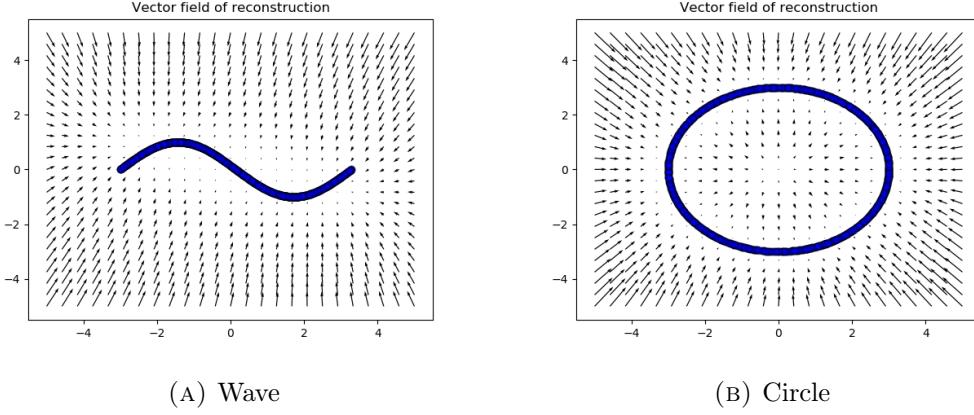


FIGURE 4.6: Vector fields of the reconstruction error evaluated on a mesh grid

The energy function and reconstruction error are computed and plotted onto heatmaps (see Figure 4.7). We can see that the two estimators have low values in the neighbourhood of the manifold and are high everywhere else. However, the reconstruction norm has also low values at the origin, which can be explained by the fact that the norm of the vectors is small at the source.

4.4 Limitations

Many interesting data structures are difficult to reproduce with shallow denoising autoencoders. For example, sequential data (e.g. sound) is better modelled with LSTM-DAE. Likewise, CNN-DAE are more appropriate to model spatial structures, such as images. However, the integrability criterion for these models is not satisfied anymore,

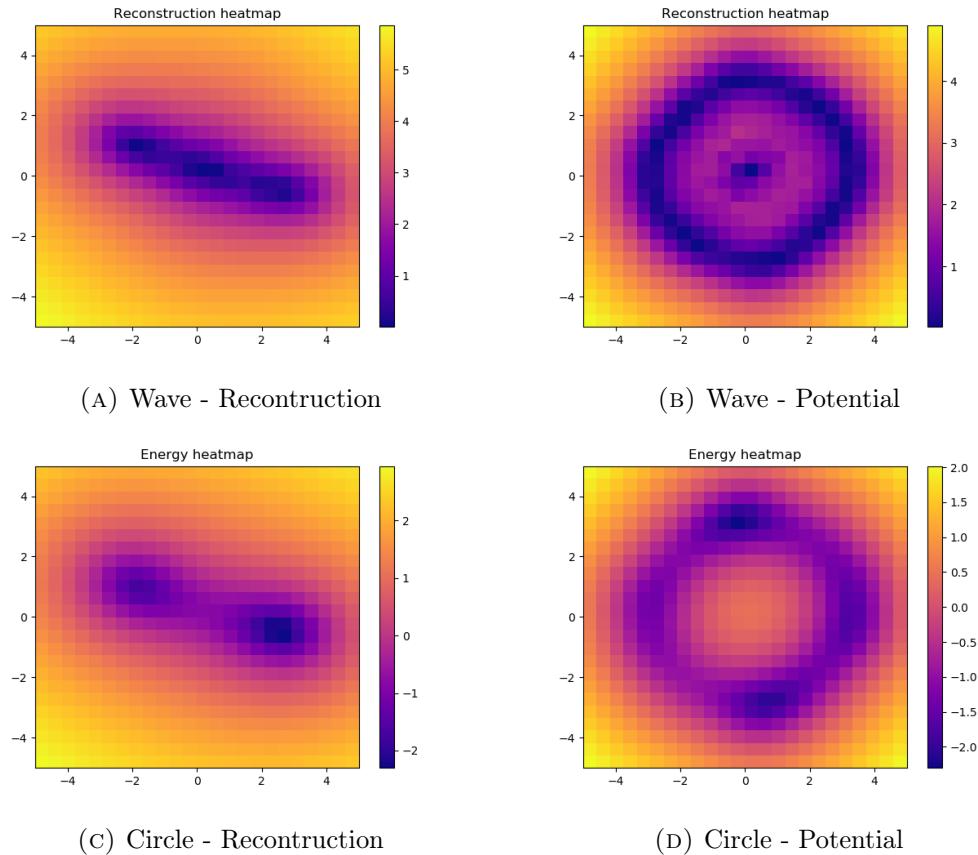


FIGURE 4.7: Estimators evaluated on wave (left) and circle (right) manifolds

and thus can not estimate the negative log-likelihood. Alternative methods to efficiently learn energy functions for spatial or sequential data are: (Zhai et al., 2016; Kim and Bengio, 2016)

Chapter 5

Energy-based Multi-Modal Attention

The literature review (Chapter 3) showed that previous research in MMDL has been mostly focused on leveraging the multi-modality to improve the accuracy of the predictions. In this chapter, a new attention module is presented to increase the robustness against failing modes: as long as at least one modality provides enough information for the task, the prediction network will be able to perform well. First, we start by providing a conceptual general framework. Following this, the design of each step of the framework is described. Finally, the training of EMMA is discussed, along with two novel regularizers.

5.1 General Framework

We define the i.i.d. dataset $\mathcal{D}^{(N)}$ with N samples (\mathbf{X}, y) . The input \mathbf{X} is composed of M modes $\{\mathbf{x}_1, \dots, \mathbf{x}_M\}$ of possibly different dimensions, such as images and sounds. The multi-modal network will be abbreviated as MMN. This model tries to make predictions \hat{y} as close as possible to the groundtruth y . The internal architecture of the MMN is often structured as a many-to-one encoder-decoder as discussed in Section 2.2. Nonetheless, the EMMA module is not constrained to any specific internal MMN architecture.

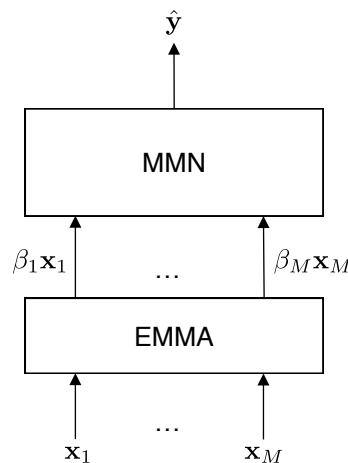


FIGURE 5.1: High-level view of a Multi-Modal Network with the EMMA module

The aim of the attention module is to spontaneously learn to assess how valuable each mode is, on a per-sample basis. The three properties influencing the importance are: relevance, failure intensity and coupling (described in Section 1.2). To learn this, modal energies are introduced, which are parametric functions embedding the three properties. The modal energy E_i of mode i is of the following form

$$E_i = f(\Psi_i) + \sum_{k \neq i}^M g[f(\Psi_i), f(\Psi_k)] \quad (5.1)$$

and will have low values if the mode of the sample is important and high values otherwise. The function f is able to capture the relationship between the relevance and the outlyingness (i.e. failure intensity), since f is optimized with respect to the loss on the predictions and is a function of the potential energy. Whereas, the role of the function g is to learn the optimal coupling between modes. Modal energies are then normalized to the importance scores via the Boltzmann distribution, thus going from a measure of absolute importance to one of relative importance. From these importance score, we can determine the attention scores β_i , representing the quantity of information that can pass through the attention module. Each mode is then multiplied by its respective attention score (see Figure 5.1). It will be made clear later on why the modes are not multiplied by the importance scores instead. Notice that the case of a missing sample in a mode is implicitly solved by definition, if we consider that missing modes are inputs with zeros.

There are two ways of interpreting the proposed solution in this chapter. First, EMMA can be seen as a sort of gate filtering perturbations out. Indeed, failing modes can provoke high activations in the MMN, disturbing the predictions. But by masking the outlying modes we diminish those activations, making it easier for the MMN to make good predictions. Another way to view it, is to understand that the MMN model easily extracts β_i and \mathbf{x}_i from the multiplication $\beta_i \mathbf{x}_i$. The model can then learn to make more robust predictions based on the extra inputs β_i .

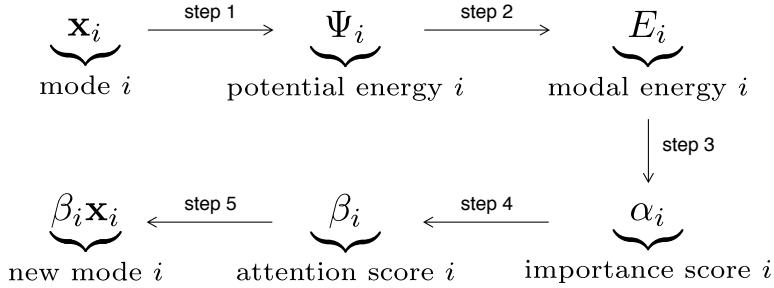


FIGURE 5.2: Summary of main steps in EMMA (step 2, 3 and 4 are detailed in the following sections, step 1 was explained in Chapter 4)

5.2 From Potential to Modal energies (step 2)

The *modal energy* of a mode is the sum of its self-energy (e_i) and the shared energies (e_{ij}) with all the other modes:

$$E_i = e_i + \sum_{k \neq i}^M e_{ij} \quad (5.2)$$

We compute the *self-energy* as a function of the potential energy,

$$e_i = w_i \Psi_i + b_i \quad (= f(\Psi_i)), \quad w_i, b_i \in \mathbb{R}^+ \quad (5.3)$$

where the parameters w_i and b_i are trained via a loss function on the predictions, in consequence the model is able via the self-energy to capture both the relevance and failure. The second advantage of this transformation is that it helps the module to face potentials on different scales, since Equation (4.7) only guarantees being proportional to the NLL, consequently potentials of different modes may not be on commensurate scales. The reason the parameters are constraint to be positive will be justified below.

Once the self-energies obtained, we can now compute the *shared energies*. The expression e_{ij} denotes the shared energy of mode j on i and is constructed from the self-energies as follows

$$e_{ij} = w_{ij} e_i^{\gamma_{ij}} e_j^{1-\gamma_{ij}} \quad (= g[f(\Psi_i), f(\Psi_j)]), \quad w_{ij} \in \left[-\frac{1}{M-1}, +\frac{1}{M-1} \right], \quad \gamma_{ij} \in [0, 1] \quad (5.4)$$

In an attempt to keep the model interpretable, we add the constraint that $\gamma_{ij} = \gamma_{ji}$. Through the use of shared energies, the model can discover the optimal coupling between the modes. If the model learns a γ_{ij} close to zero, mode i and j will influence each other much more than a γ_{ij} near to one. In other words, the parameter γ_{ij} learns the degree of coupling in the spectrum from strongly coupled ($\gamma_{ij} = 0$) to independent ($\gamma_{ij} = 1$). The direction of coupling between mode i and j are learned by the weights w_{ij} and w_{ji} . For a positive w_{ij} , an increase in self-energy e_j causes an increase in e_{ij} , and thus an increase in e_i . Whereas if w_{ij} is negative, an increase in e_j leads to a decrease in e_{ij} . The weights w_{ij} are not imposed to be equal to w_{ji} , such that modes can influence each other asymmetrically. This assymmetry is justified by the following example: take a multi-modal problem with three modes A, B and C. We want the model to learn that if mode A is failing, it is optimal that mode B "takes over". And if mode B is failing, it is optimal for C to "take over". This example can only be modelled with asymmetry.

A consequence of the design of Equation (5.4) is that the evaluation of the gradient during the backpropagation step now involves taking the logarithm of e_i ¹, which is undefined for negative values. As the weights in Equation (5.3) are positive, we only have to make sure the values of the potential energy are positive. The latter is done by lowering the potential Ψ_i to Euler's number e as

$$\Psi_i \leftarrow \max(e, \Psi_i - \Psi_i^{(\min)} + e) \quad (5.5)$$

where $\Psi_i^{(\min)}$ denotes the lowest value of Ψ_i in the training set. This correction avoids undefined values ($\Psi_i \geq 0$) but also exploding gradient ($\Psi_i \geq e$). The reason a max-operator is used is because lower energy values than $\Psi_i^{(\min)}$ can occur during inference. Of course this correction must be performed prior to the computation of self-energies.

¹See Appendix B.2

5.3 From Modal energies to Importance scores (step 3)

The importance scores are computed from the modal energies via the Boltzmann distribution:

$$\alpha_i = \frac{1}{Z} e^{-\rho E_i} \quad \text{with the partition function} \quad Z = \sum_{k=1}^M e^{-\rho E_k} \quad (5.6)$$

This guarantees the scores to be normalized and summing up to one. A mode i will be said to be important if its score is close to one (low modal energy E_i). The hyperparameter ρ represents the coldness, the inverse of the temperature. It controls the entropy of the importance scores distribution. At high temperature ($\rho \rightarrow 0$) the distribution becomes more uniform, and at low temperature ($\rho \rightarrow +\infty$) the importance scores corresponding to the lowest energy tends to 1, while the others approach 0. As can be observed on Figure 5.3, the coldness has a significant influence on the overall behaviour of the attention module; Careful tuning of ρ is thus necessary.

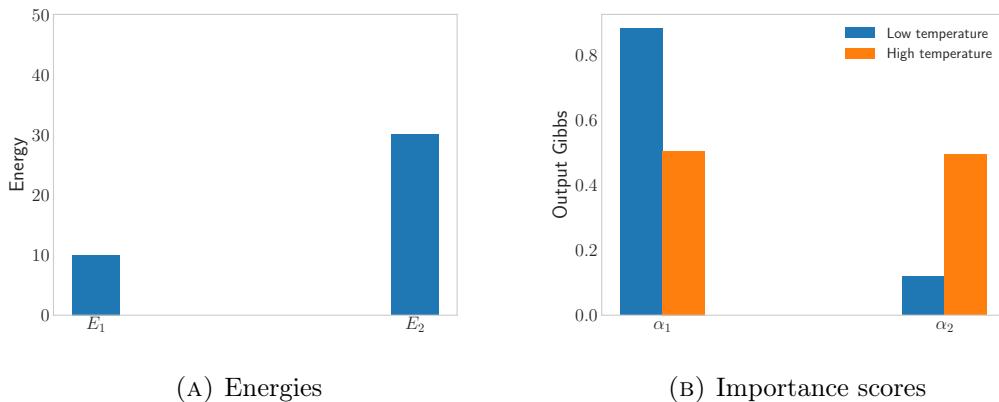


FIGURE 5.3: Input-output of Boltzmann distribution for two different temperatures, low temperature ($\rho = 0.1$) and high temperature ($\rho = 0.001$)

5.4 From Importance to Attention scores (step 4)

The attention scores are given by

$$\beta_i = \tanh(g_a \alpha_i - b_a) \quad \text{with} \quad g_a > 0, \quad b_a \in [0, 1] \quad (5.7)$$

The hyperbolic tangent adds non-linearity while the gain g_a and bias b_a enable the model to control the threshold and capacity (see Figure 5.4). The latter two concepts are detailed below.

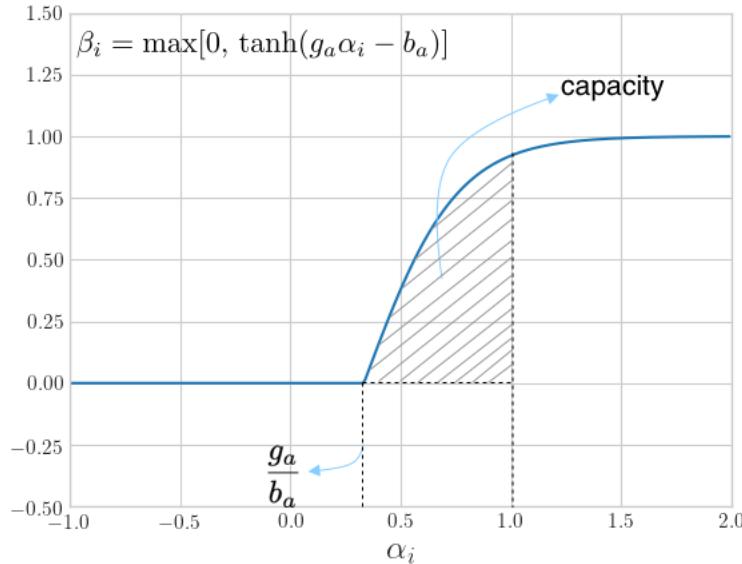


FIGURE 5.4: Attention function (the max-operator generalizes the attention function to cases where $\alpha \in \mathbb{R}$)

Energy threshold

The module will let the information of mode i pass by, only if $g_a \alpha_i - b_a > 0$

$$\begin{aligned} &\Leftrightarrow \log(\alpha_i) > \log(b_a/g_a) \\ &\Leftrightarrow E_i \geq \frac{\log(g_a/b_a) - \log(Z)}{\rho} = E_{\text{threshold}} \end{aligned} \quad (5.8)$$

where $E_{\text{threshold}}$ represents the maximal amount of energy is allowed to have, to not be completely masked off (see Figure 5.4). We deduce that the learned gain and bias control this threshold. Nevertheless, the value of the partition function Z also controls the threshold, making it dynamic. The partition function will be higher if the total energy ($\sum_i E_i$) is higher, resulting in a diminishment of the threshold. To put it in another way, EMMA adapts the selectiveness with respect to the overall quality of the entire input sample. Notice that the influence of the temperature (ρ^{-1}) is non-trivial to analyse, because Z also depends on ρ .

Capacity

A more common way to write the attention function would be $\tanh(\mathbf{W}\boldsymbol{\alpha} + \mathbf{b})$, whereas we have $\tanh(g_a \mathbf{I}\boldsymbol{\alpha} - b_a \mathbf{u})$ with the unit vector $\mathbf{u} = (1 \dots 1)^T$. We argue the latter better mimics human's attention, permitting us to introduce the concept of capacity, which in psychology is viewed as the amount of resource that can be allocated (Kahneman, 1975). If we look at Figure 5.4, this can be translated as,

$$\text{capacity} \triangleq \int_0^1 \tanh(g_a \alpha - b_a) d\alpha \quad (5.9)$$

Define the auxiliary variable $u = g_a\alpha - b_a$. Now using

$$\frac{du}{d\alpha} = g_a \Leftrightarrow d\alpha = \frac{1}{g_a} du \quad (5.10)$$

we can write

$$\begin{aligned} \text{capacity} &= \frac{1}{g_a} \int_0^1 \tanh(u) du \\ &= \frac{1}{g_a} \log[\cosh(g_a\alpha - b_a)] \Big|_{\alpha=0}^1 + \text{constant} \\ &= \frac{1}{g_a} \log \left[\frac{\cosh(g_a - b_a)}{\cosh(-b_a)} \right] \end{aligned} \quad (5.11)$$

When the capacity is too low, no sufficient amount of information is passed to the MMN, leading to wrong predictions. Similarly, if the capacity is too high, the perturbations of the failing modes will pass and cause a decrease in performances. It is expected that the model learns the optimal trade-off, however, if we want the attention module to be robust against failing situations it was not trained on, it can be interesting to try to control this trade-off. To this end, we created a simple regularizer which is discussed in the next section. Observe that the concept of capacity can also be applied to $\tanh(\mathbf{W}\alpha + \mathbf{b})$, but each mode would have his own capacity, making the importance scores less meaningful.

5.5 Training & Regularization

The training of the attention module and the prediction model is performed in two stages (see Figure 5.5). First, each mode is assigned a separate autoencoder, which is trained on the mode to learn the potential energy function. Once trained, the weights of the autoencoders are freezed. In the second phase, EMMA is inserted in front of the MMN and is trained end-to-end on both normal and corrupted data. By corrupted data, we mean samples on which a corruption process is applied in order to simulate one or more failing modes.

Additionaly, two regularizers are introduced. The first one controls the capacity, where λ_c can be positive/negative depending on if we want to maximize/minimize the capacity.

$$\tilde{\mathcal{L}} = \mathcal{L}(y, \hat{y}) + \lambda_c g_a - \lambda_e \Omega \quad \text{with} \quad \Omega = \sum_{k=1}^M \xi_k \log(\alpha_k) \quad \text{and} \quad \xi_k = \begin{cases} \xi_- = -1 & \text{if } \mathbf{x}_k \text{ is corrupted} \\ \xi_+ = +1 & \text{otherwise} \end{cases} \quad (5.12)$$

Secondly, the purpose of regularizing the energy ($\lambda_e \Omega$) is to control the trade-off between, on one side the relevance and coupling, and on the other side the failure intensity. Without a regularizer, the parameters of the modal energy functions are optimized only regarding the predictions (\mathcal{L}), which could lead to a large discrepancy between modal energies E_i and their original potential energies Ψ_i . Although the energy regularizer is relatively straightforward, we will show below that some care needs to be taken regarding the corruption process.

Energy regularization

Let $\boldsymbol{\theta}$ be the set of all the parameters in step 2 of the attention module. The effect of the energy regularizer in the SGD algorithm is isolated and written

$$\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} + \epsilon \lambda_e \nabla_{\boldsymbol{\theta}} \Omega \quad (5.13)$$

Remember the objective, we want this update to minimize the discrepancy, thus decrease/increase modal energies E_i for low/high potential energies Ψ_i . To verify this let us compute² $\nabla_{\boldsymbol{\theta}} \Omega$,

$$\nabla_{\boldsymbol{\theta}} \Omega = \sum_{k=1}^M \xi_k \nabla_{\boldsymbol{\theta}} \log(\alpha_k) \quad (5.14)$$

The gradient of the logarithm can be developed as

$$\begin{aligned} \nabla_{\boldsymbol{\theta}} \log(\alpha_k) &= \nabla_{\boldsymbol{\theta}} \log \left[\frac{e^{-\rho E_k}}{Z} \right] \\ &= \nabla_{\boldsymbol{\theta}}(-\rho E_k) - \nabla_{\boldsymbol{\theta}} \log \sum_{l=1}^M e^{-\rho E_l} \\ &= -\rho \nabla_{\boldsymbol{\theta}} E_k - \frac{\sum_{l=1}^M \nabla_{\boldsymbol{\theta}} e^{-\rho E_l}}{\sum_{l=1}^M e^{-\rho E_l}} \\ &= -\rho \nabla_{\boldsymbol{\theta}} E_k + \rho \frac{\sum_{l=1}^M e^{-\rho E_l} \nabla_{\boldsymbol{\theta}} E_l}{\sum_{l=1}^M e^{-\rho E_l}} \\ &= \rho \left[-\left(1 - \frac{e^{-\rho E_k}}{Z}\right) \nabla_{\boldsymbol{\theta}} E_k + \sum_{l \neq k}^M \frac{e^{-\rho E_l}}{Z} \nabla_{\boldsymbol{\theta}} E_l \right] \\ &= \rho \left[-(1 - \alpha_k) \nabla_{\boldsymbol{\theta}} E_k + \sum_{l \neq k}^M \alpha_l \nabla_{\boldsymbol{\theta}} E_l \right] \end{aligned} \quad (5.15)$$

We go further by expressing the equation above with respect to the subset of parameters $\boldsymbol{\theta}_i = \{[\gamma_{ik}, w_{ik}]_{k=1}^M, w_i, b_i\}$:

$$\nabla_{\boldsymbol{\theta}_i} \log(\alpha_k) = \begin{cases} -\rho(1 - \alpha_i) \nabla_{\boldsymbol{\theta}_i} E_i, & \text{if } i = k \\ \rho \alpha_i \nabla_{\boldsymbol{\theta}_i} E_i, & \text{if } i \neq k \end{cases} \quad (5.16)$$

The gradient of the regularizer can now be computed by plugging Equation (5.16) into the summation (5.14). Let M' be the number of uncorrupted modes. We obtain for an uncorrupted mode i ,

$$\nabla_{\boldsymbol{\theta}_i} \Omega = \xi_+ \left[-\rho(1 - \alpha_i) \nabla_{\boldsymbol{\theta}_i} E_i \right] + \left[(M' - 1)\xi_+ + (M - M')\xi_- \right] \alpha_i \rho \nabla_{\boldsymbol{\theta}_i} E_i \quad (5.17)$$

and for a corrupted mode i ,

$$\nabla_{\boldsymbol{\theta}_i} \Omega = \xi_- \left[-\rho(1 - \alpha_i) \nabla_{\boldsymbol{\theta}_i} E_i \right] + \left[M'\xi_+ + (M - M' - 1)\xi_- \right] \alpha_i \rho \nabla_{\boldsymbol{\theta}_i} E_i \quad (5.18)$$

²the batch is assumed to only contain one sample for the sake of simplicity. However, the demonstration can be generalized to any batch size.

Substituting ξ_k , we can summarize Equations (5.17) and (5.18) as

$$\nabla_{\theta_i} \Omega = -[(M - 2M')\alpha_i + \xi_i]\rho \nabla_{\theta_i} E_i \quad (5.19)$$

Adding the constraint that $M' = \lfloor \frac{M+1}{2} \rfloor$, two cases can be distinguished. If the total number of modes M is even, then we have

$$\theta_i \leftarrow \theta_i - \epsilon \lambda_e \rho \xi_i \nabla_{\theta_i} E_i \quad \text{with } \lambda_e \in \mathbb{R}^+ \quad (5.20)$$

Ignoring the second-order effects of the Taylor expansion of the modal energy function, it can be concluded from the equation above that the regularizer will update the parameters such that the values of the modal energy E_i increases/decreases if mode i is corrupted/uncorrupted.

In analogy, if M is uneven we have

$$\theta_i \leftarrow \begin{cases} \theta_i - \epsilon \lambda_e \rho (1 - \alpha_i) \nabla_{\theta_i} E_i, & \text{if } i \text{ is uncorrupted} \\ \theta_i + \epsilon \lambda_e \rho (1 + \alpha_i) \nabla_{\theta_i} E_i & \text{otherwise} \end{cases} \quad (5.21)$$

The principle is the same as in the even case with an additional effect: the correction will be proportional to the error. To put it in another way, high energies that must be low and low energies that have to be high will have stronger gradients than their counterparts. This is similar to the positive and negative phase in the optimization of Restricted Boltzmann Machines.

To conclude, let us notice that some undesired effects can appear if we do not add the constraint $M' = \lfloor \frac{M+1}{2} \rfloor$. As an illustration, take $M' = \lfloor \frac{M+1}{2} \rfloor + 1$, Equation (5.13) becomes

$$\theta_i \leftarrow \theta_i - \epsilon \lambda_e \rho (\alpha_i + \xi_i) \nabla_{\theta_i} E_i \quad (5.22)$$

which is unstable for uncorrupted modes leading to a collapse where all energies tend to decrease.

5.6 Advantages

The key advantages of using EMMA are:

- The generic design of EMMA permits it to be easily added to any type of architecture of a multi-modal model, without modifying nor EMMA nor the MMN.
- The burden on the MMN is reduced, it only has to learn to make good predictions from the received information. The MMN does not need anymore to learn to distinguish failing modes.
- The interpretability is increased; It can be verified on a per-sample basis which modes are used to make predictions using the attention scores. As we will see in Chapter 6, EMMA gives us a measure ($\sum_i E_i$) of how uncertain the MMN is about its predictions using the total energy.

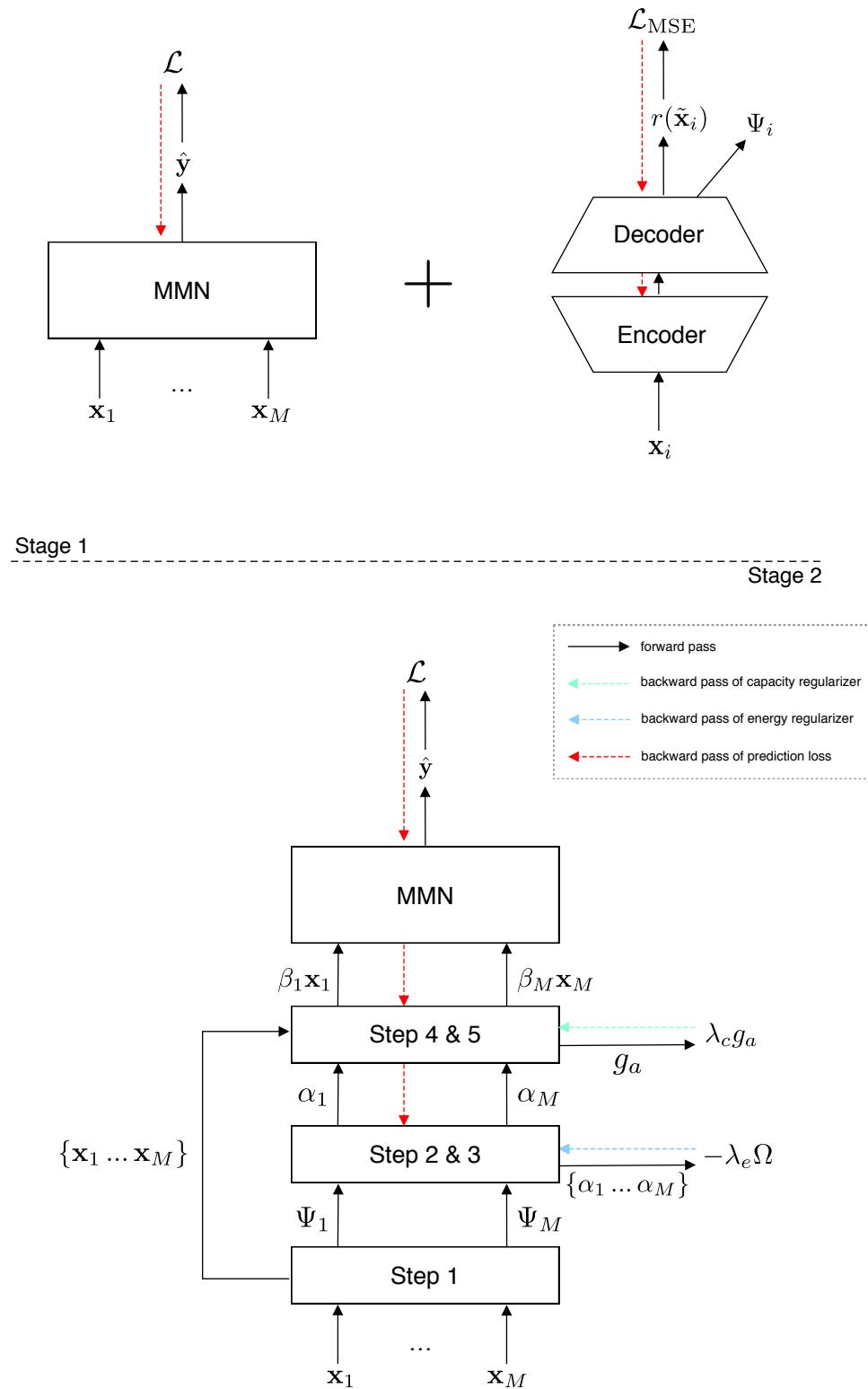


FIGURE 5.5: Summary of end-to-end training

Chapter 6

Experiments & Results

Both are experiments on dataset described in previous chapter. Experiment II will be about energy estimation. Experiment III evaluates and analyzes the robustness of the model with and without EMMA.

6.1 Pulsar detection

The models will be trained to detect pulsar stars. In (very) short, pulsar stars are neutron stars emitting radio waves on a periodic time-frame. A summary of the seminal work of (Lyon, 2016) can be found in Appendix A. The thesis can be accessed [here¹](#) and the dataset [here²](#).

Short description of two modes (ip and dm) and internal features (mean...). Explain why detection is difficult. Classification signal/background. Skewed dataset, give numbers.

6.2 Corruption

- standardize, why? split sets and apply one from train [error standardize](#)
- SNR (see good explanation in pulsar thesis). If greater than 1, signal non-distinguishable. White noise. Explain it is not the same than AE corruption. $10 \log(\frac{1}{\sigma^2})$
- on signal and background because we corrupt the whole mode and not the class

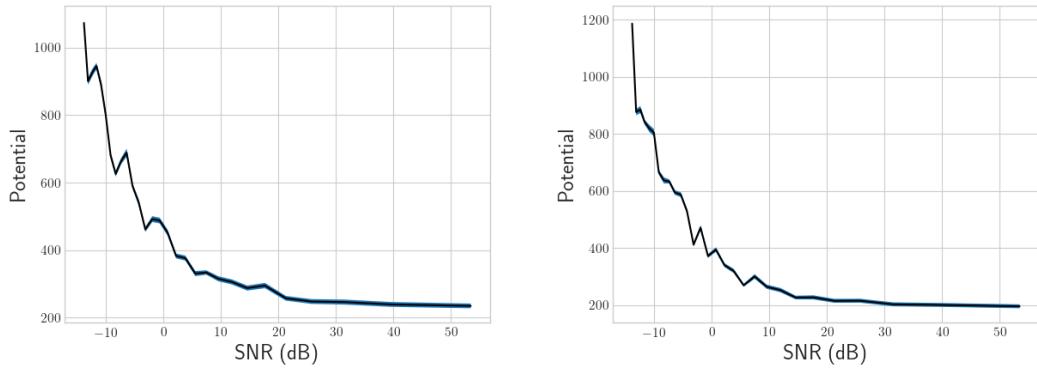
6.3 Experiment II

- train-test split
- AE trained on train-set and then test on test-set
- matrix with number of signals, ... (eda)

Setup: max epochs = 30, batch size = 64, noise DAE = 0.01, d input = 4, n hidden = 12, adam 0.001, sigmoid

¹http://www.scienceguyrob.com/wp-content/uploads/2016/12/WhyArePulsarsHardToFind_Lyon_2016.pdf

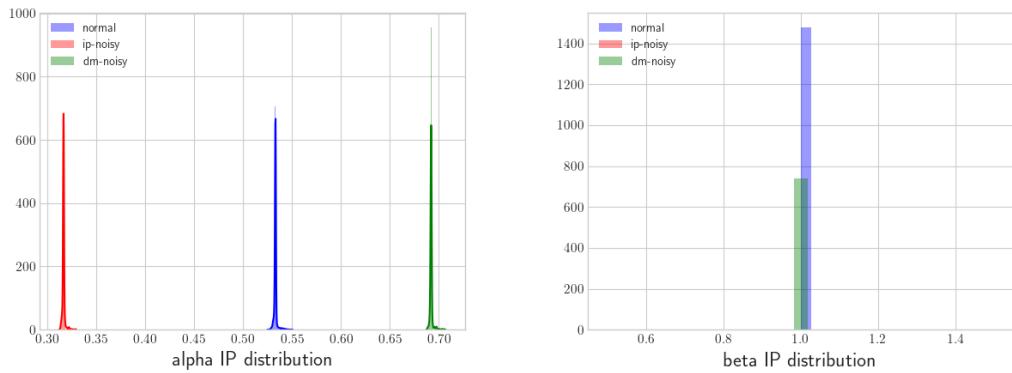
²<https://archive.ics.uci.edu/ml/datasets/HTRU2>



6.4 Experiment III

- BCE & F1 (not AUC – explain why). [F1vsAUC1](#). [F1vsAUC2](#). [F1vsAUC3](#). [F1](#)
- train-valid-test
- threshold optimal choice via ROC. all on valid set
- 3 models: base (train normal, valid normal), without (train noisy, valid noisy), with (train noisy, valid noisy)
- 50-25-25 noisy mode. give detailed numbers. eda.
- trained with early stopping + retrain for .. epochs with valid+train. saved model.
- one subsection per plot: explain details experiment and how results are obtained. then analyze and conclusions.

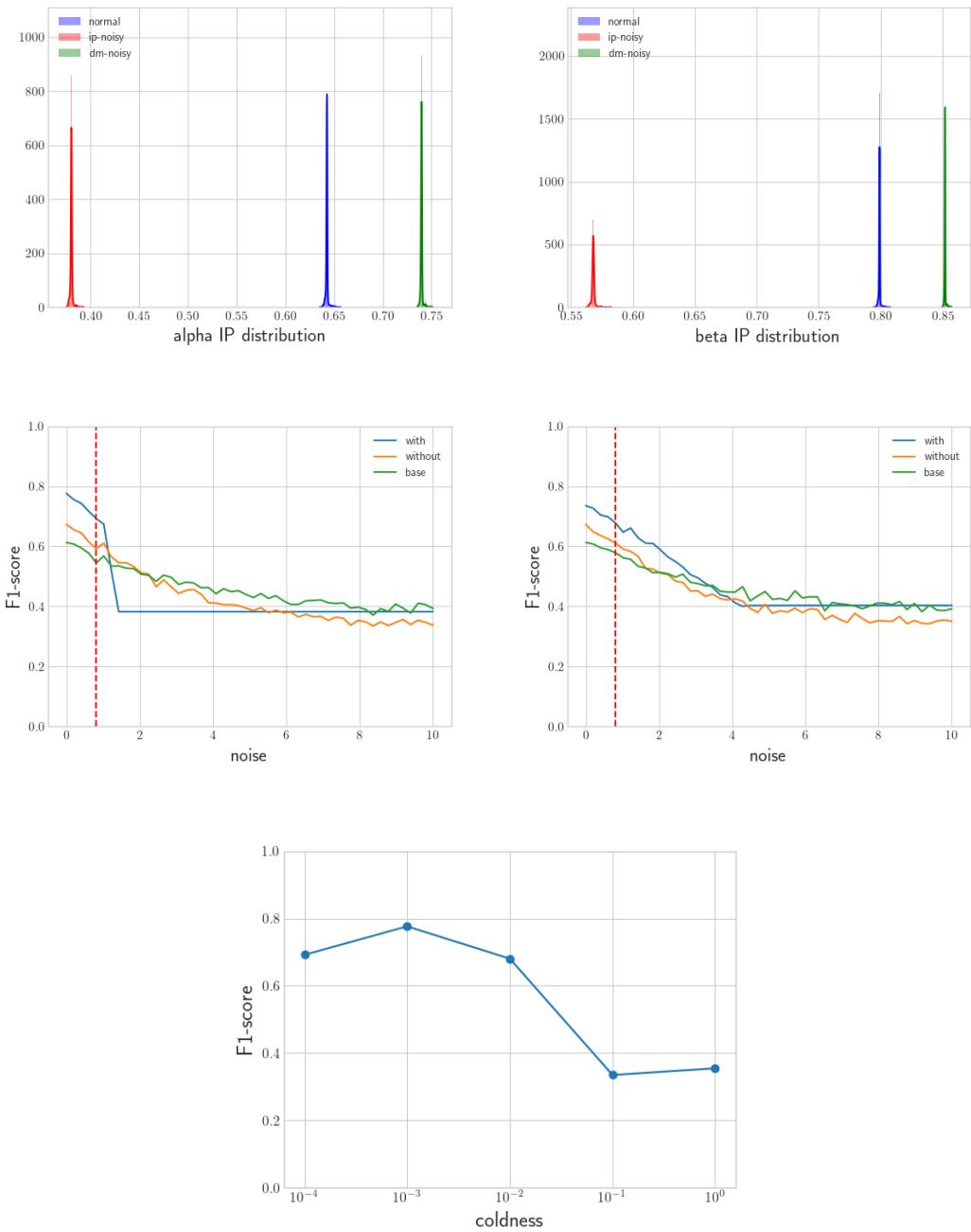
Attention-shift



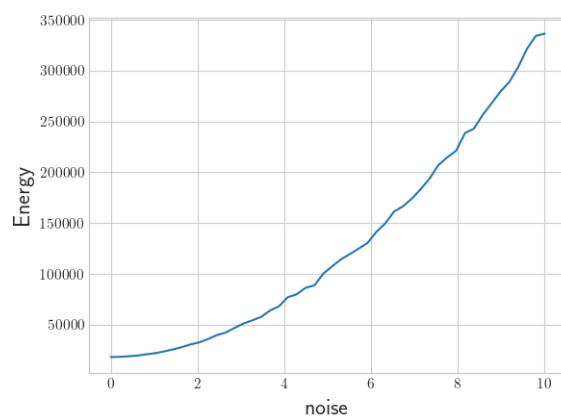
Robustness generalisation

Yerkes-Dodson curve

over-under arousal. do on larger range.



Energy generalisation



Chapter 7

A Unified Model for Multi-Modal Attention

The purpose of using EMMA is to help the multi-modal network (MMN) to handle failing modes, but also as a side effect, modes with different contributions and modes with different levels of noise. Chapter 3 discussed self-attention and crossmodal attention that are used to highlight information inside a specific mode, such as certain regions in an image or a set of frequencies in a sound. The difference between the two is that self-attention uses only the information of the mode itself as a context, whereas crossmodal attention leverages the information in all the modes. We claim to have all the ingredients to construct a complete multi-modal network. As a reminder, human's multi-modal attention consists of three different components: exogenous, endogenous and crossmodal attention. Attention is endogenous when we voluntary choose to attend to something whereas exogenous occurs when a person's attention is captured reflexively by the sudden onset of an unexpected event (Driver and Spence, 1998). Thus, an endogenous module could easily be constructed as a block of M self-attentions. Moreover, EMMA can be considered to be equivalent to exogenous attention.

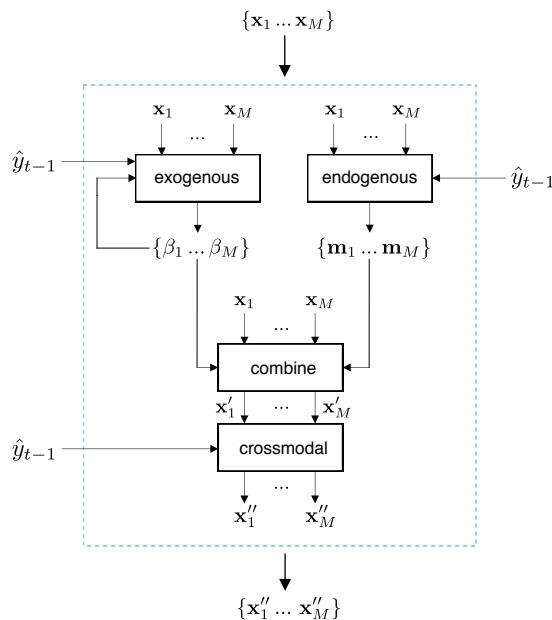


FIGURE 7.1: A possible architecture for a unified multi-modal attention

With this in mind, we present a unified model (see Figure 7.1) combining all the strengths of each type of attention. First, the attention masks of the exogenous model,

β_i , and the attention masks of the endogenous module, \mathbf{m}_i , are combined as $\beta_i \mathbf{m}_i$. The resulting mask is then applied to the input sample, and is passed through the crossmodal module, which explores the relationships between the modes of the input. Finally, the processed input $\{\mathbf{x}'_1 \dots \mathbf{x}'_M\}$ is forwarded to the MMN. In addition, the complete module can be further refined by inserting feedback loops from the predicted output to the separate modules as it is often done in the literature (Afouras et al., 2018; Vaswani et al., 2017; Bahdanau, Cho, and Bengio, 2014). For example, in a self-driving system, the model could adapt its focus to different regions of the input image depending on the previously detected cars. To conclude, let us emphasize that the proposed architecture is only a generalization of attention networks such as (Afouras et al., 2018), supplemented with an exogenous component.

Chapter 8

Conclusion

Summary of what was seen/done during the master thesis from start to end. Scales up quadratically with number of modes, explain.

8.1 Contributions

Summarize contributions

8.2 Research questions

Each research question + answer

8.3 Future work

- Annealing + init, end temperature + explain init of next layer [Linke annealing](#). Multiple modes then blows up. Image 100 modes (although not realistic in real-world problems) then tend to zero. Solution add a common gain? Analyze influence of multiple modes etc..
- Explore different shared energies design
- Images/sound sequences. early, late fusion, manifold with respect to unified model? it could be very easy to test on images/sounds (give even during test-time a true outlyingness measure, not specifically the NLL) and at the same time investigate general ways of getting such a measure.
- Unseen values

Appendix A

Dataset

This part of the thesis provides a brief overview of what pulsar stars are. It then goes on to explain the two modes of the dataset used for their detection: the integrated profile (IP) and dispersion measure (DM). The most part of this chapter is a direct summary of the background chapter in the doctoral thesis (Lyon, 2016). Some parts of the text are barely changed from (Lyon, 2016), nevertheless, many details have been voluntary ignored for the sake of simplicity as it is not the focus of this work.

A star is a luminous ball of gas, mostly hydrogen and helium, held together by its own gravity. The nearest star of Earth is the Sun. Most mass is in the core, at the center of the star. Gravitational forces are by consequence directed inwards. During the majority of a star's life, it will fuse hydrogen to helium, generating an outwards pressure, balancing the gravity (Ghosh, 2007). By the time the hydrogen amounts become insufficient, the star starts to use other elements in the surrounding layers of the core as a fuel. As those elements diminish, the star's energy output drops rapidly, causing gravity to overcome the forces which had previously maintained the stars structure. The core of the star than undergoes a rapid and violent collapse (Ghosh, 2007). The collapse can lead to a number of potential evolutionary outcomes for the leftover core (see Figure A.1), depending on the stars birth mass measured in solar masses (M_{\odot}). Intuitively, the heavier the birth mass, the greater the inwards gravitational force are and the harder the collapse. The first outcome applies to low mass stars, which typically become white dwarfs following their collapse. Within white dwarfs, densely packed electrons are able to resist gravitational compression. Our own sun is likely to one day become a white dwarf star. Then there are stars between 8-20 M_{\odot} at birth, electron degeneracy pressure can no longer prevent collapse as in white dwarfs, but they are not massive enough to undergo complete gravitational collapse, preventing the formation of a black hole. Instead the intense conditions within these stars cause electrons to combine with protons forming neutrons who resist against pressure; These stars are called neutron stars. The last evolutionary outcome applies to large stars with masses greater than 20 M_{\odot} . These stars can, under the right conditions, undergo complete gravitational collapse. This results in the formation of a black hole singularity otherwise known as a stellar mass black hole.

A pulsar is a unique form of neutron star that retained most of its angular momentum of their progenitor star during collapse. Complex interactions between the surfaces of pulsars and their strong magnetic fields, helps to produce their defining feature, the emission of radio waves. The radio emission produced by pulsars originates from their magnetospheres (Ghosh, 2007). This is the area of space surrounding a pulsar in which charged particles are influenced by a co-rotating magnetic field, which has both open and closed field lines (D.R. and M., 2005) (see Figure A.2). To maintain this co-rotation

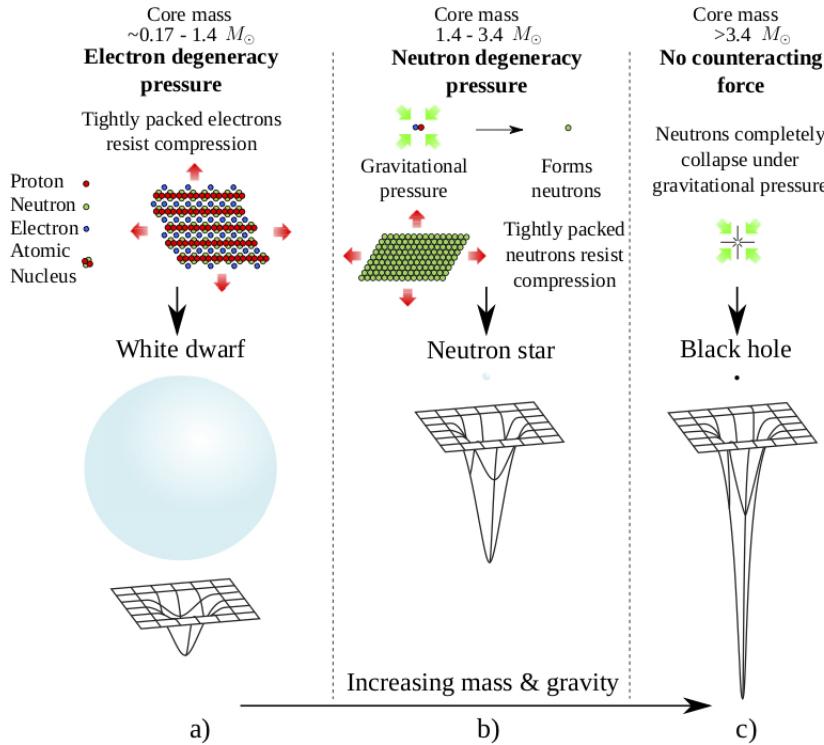


FIGURE A.1: Common evolutionary endpoints for main sequence stars. In a) electron degeneracy pressure prevents gravitational collapse, leading to the formation of a white dwarf star. In b) electron degeneracy pressure is no longer enough to counteract the inward force of gravity, however the gravitational pressure is insufficient to overcome neutron degeneracy pressure, allowing a Neutron star to form. Finally in c) the force of gravity is so great that gravitational collapse cannot be halted, resulting in the formation of a black hole. The depictions of the gravitational sinks above are based on diagrams by (Treat and Stegmaier, 2014). *Image and caption copied from* (Lyon, 2016).

property, the velocity of the field lines must increase as they move further away from the pulsar. Eventually the distance becomes so great, that to maintain co-rotation, the velocity of the field lines must be greater than or equal to the speed of light c . This is not possible, thus the field lines are unable to close where the required velocity is c . The abstract cylinder aligned with the rotation axis, that synchronously rotates with the pulsar at a velocity c , is known as the light cylinder (see Figure A.2). The particles extracted from the surface are then believed to be accelerated along the co-rotating magnetic field lines of the magnetosphere (Lorimer, 2008), which endows the particles with increased energy. This additional energy causes the particles to emit radiation (Lorimer, 2008) to be emitted along the open field lines near a pulsar's magnetic pole. A pulsar's magnetic axis is usually inclined with respect to its rotational axis. Therefore each time a pulsar rotates, the radiation beam produced near the magnetic poles, is swept at an angle across the sky. If the beam crosses the line of sight of an observer here on Earth, the pulsar becomes detectable as a rise and fall in broadband radio emission. This pattern repeats periodically with each rotation of the pulsar. This is known as the lighthouse model of emission (Lorimer, 2008), because the beam of radiation is analogous to a lighthouse warning light rotating very quickly.

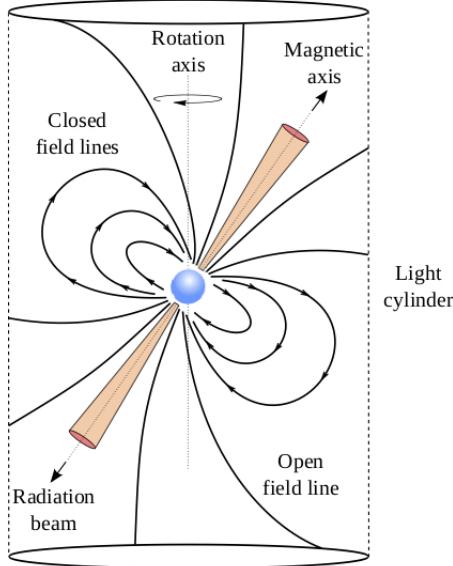


FIGURE A.2: Simplification of the lighthouse model of a radio pulsar, the pulsar is surrounded by a strong magnetic field comprising of open and closed field lines unable to close at the light cylinder. The light cylinder is an imaginary cylinder aligned with the pulsar's rotational axis, that synchronously rotates with the pulsar at the speed of light. As the magnetic field cannot rotate at this velocity, the field lines cannot close at the light cylinder leading to open field lines. Radio pulses are emitted from the open field lines at a region near the magnetic poles in the pulsar's magnetosphere. *Image and caption copied from (Lyon, 2016).*

Each pulsar produces a unique pattern of pulse emission known as its pulse profile (Lorimer, 2008). Two such profiles are shown in Figure A.3. However whilst pulsar rotational periods are extremely consistent, their profiles can deviate from one-period to the next. Whilst such changes in the pulse profile provide clues to what is happening in and around the pulsar, they make pulsars hard to detect. This is because their signals are non-uniform and not entirely stable overtime. However these profiles do become stable, when averaged over many thousands of rotations.

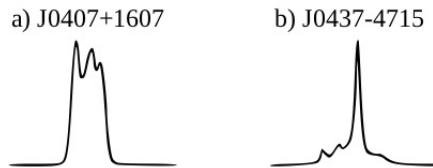


FIGURE A.3: Example pulse profiles of two separate pulsars. These profiles were adapted from those originally presented in (D.R. and M., 2005). *Image and caption copied from (Lyon, 2016).*

Signals travelling through the interstellar medium (ISM) are affected, the most significant effect is known as dispersion. As pulsar signals travel through the ISM towards the Earth, they interact with charged particles (free electrons) on route. These interactions delay the arrival of the signal here on Earth. The low frequency components of the signal are delayed more than the corresponding high frequency counterparts. This has a dispersive effect that causes pulsar signals to become smeared in time. This makes it

difficult to detect pulsars, as their pulses become less pronounced as shown in Figure A.4. Manifesting itself as a reduction in the signal-to-noise ratio of a detected pulse. The amount of dispersive smearing a signal receives is proportional to a quantity called the dispersion measure (DM) (D.R. and M., 2005). The DM is the integrated column density of free electrons between an observer and a pulsar (Lorimer, 2008). The true column density, and thus the precise degree to which a signal is dispersed, cannot be known a priori. A number of dispersion measure tests or "DM trials", must therefore be conducted to determine this value as accurately as possible. An accurate DM can be used to undo the dispersive smearing, allowing the signal-to-noise ratio of a detected signal to be maximised (D.R. and M., 2005). For a single dispersion trial, each frequency channel is shifted by an appropriate delay. Subsequent trials increment the delay in steps, until a maximum DM is reached. This maximum will vary according to the region of sky being surveyed, the observing frequency, and bandwidth. The process produces one 'de-dispersed' time series per frequency channel. These are then summed to produce a single de-dispersed time series per trial (as shown at the bottom plot of Figure A.4 a). In total de-dispersion produces a number of time series equal to the total number of DM trials.

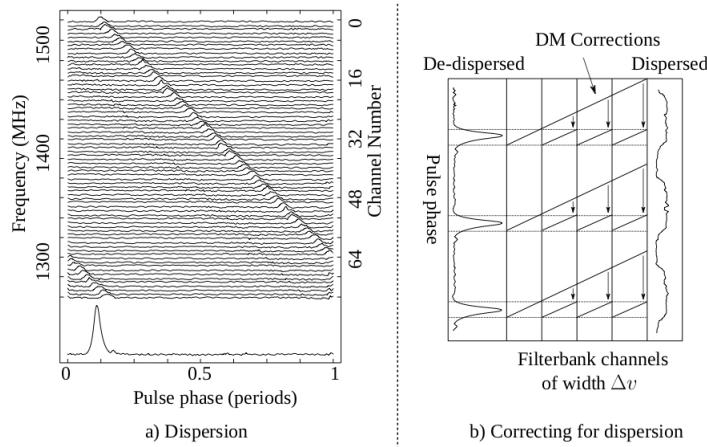


FIGURE A.4: An example of signal dispersion. Based upon diagrams originally presented in (D.R. and M., 2005). Plot a) shows how a signal is dispersed in time. Dispersion hides the true pulse shape and causes a lowering of the detected signal-to-noise. Plot b) shows the application of DM corrections to a dispersed signal. The DM correction is different in each frequency channel, since dispersion is proportional to frequency.

Image and caption copied from (Lyon, 2016).

By precisely measuring the timing of such pulses, astronomers can use pulsars for unique experiments at the frontiers of modern physics. Indeed, pulsars exist in strong-field gravitational environments due to their enormous mass. It is impossible to study such environments within Earth-based laboratories, or even within the confines of our own solar system which is lightweight by comparison. In the strong-field environment provided by pulsars, their immense gravitational fields directly affect the arrival times on Earth of the signals they produce, via special and general-relativistic effects. By studying these effects, tests of many gravitational theories can be accomplished. Another application of measuring the arrival time of pulses is that they are effective time keeping system, rivalling atomic clocks for accuracy. Such clocks are useful for spacecraft navigation and timekeeping here on Earth.

Appendix B

Miscellaneous

B.1 Integrability criterion

The integrability criterion (Santilli, 1982) is a sufficient condition for a vector field to be a gradient field as well. It states that for some open, simple connected set U , a continuously differentiable function $F : U \rightarrow R^L$ defines a gradient field if and only if

$$\frac{\partial F_j(\mathbf{x})}{\partial x_i} = \frac{\partial F_i(\mathbf{x})}{\partial x_j}, \quad \forall i, j = 1 \dots L \quad (\text{B.1})$$

In other words, integrability follows from the symmetry of the partial derivatives.

B.2 Gradient with respect to gamma

The gradient of the loss with respect to the coupling parameter γ_{ij} is computed with the chain rule:

$$\frac{\partial \tilde{\mathcal{L}}}{\partial \gamma_{ij}} = \frac{\partial \tilde{\mathcal{L}}}{\partial E_i} \cdot \frac{\partial E_i}{\partial \gamma_{ij}} + \frac{\partial \tilde{\mathcal{L}}}{\partial E_j} \cdot \frac{\partial E_j}{\partial \gamma_{ij}} \quad (\text{B.2})$$

In particular,

$$\begin{aligned} \frac{\partial E_i}{\partial \gamma_{ij}} &= \frac{\partial}{\partial \gamma_{ij}} \sum_{k=1}^M E_{ik} \\ &= \frac{\partial E_{ij}}{\partial \gamma_{ij}} + \frac{\partial E_{ji}}{\partial \gamma_{ij}} \\ &= \frac{\partial}{\partial \gamma_{ij}} (w_{ij} e_i^{\gamma_{ij}} e_j^{1-\gamma_{ij}}) + \frac{\partial}{\partial \gamma_{ij}} (w_{ji} e_j^{\gamma_{ij}} e_i^{1-\gamma_{ij}}) \\ &= w_{ij} e_i^{\gamma_{ij}} \frac{\partial}{\partial \gamma_{ij}} e_j^{1-\gamma_{ij}} + e_j^{1-\gamma_{ij}} \frac{\partial}{\partial \gamma_{ij}} e_i^{\gamma_{ij}} + \dots \\ &= (\log e_i + \log e_j) (w_{ij} e_i^{\gamma_{ij}} e_j^{1-\gamma_{ij}} + w_{ji} e_j^{\gamma_{ij}} e_i^{1-\gamma_{ij}}) \end{aligned} \quad (\text{B.3})$$

As we can see, the gradient with respect to γ_{ij} does indeed involve a natural logarithm of self-energies e_i and e_j . Thus, self-energies must be constrained to positive values.

Bibliography

- Afouras, Triantafyllos et al. (2018). “Deep Audio-Visual Speech Recognition”. In: *arXiv e-prints*, arXiv:1809.02108, arXiv:1809.02108. arXiv: [1809.02108 \[cs.CV\]](#).
- Alain, Guillaume and Yoshua Bengio (2012). “What Regularized Auto-Encoders Learn from the Data Generating Distribution”. In: *arXiv e-prints*, arXiv:1211.4246, arXiv:1211.4246. arXiv: [1211.4246 \[cs.LG\]](#).
- Bahdanau, Dzmitry, Kyunghyun Cho, and Yoshua Bengio (2014). “Neural Machine Translation by Jointly Learning to Align and Translate”. In: *arXiv e-prints*, arXiv:1409.0473, arXiv:1409.0473. arXiv: [1409.0473 \[cs.CL\]](#).
- Baltrušaitis, T., C. Ahuja, and L. Morency (2019). “Multimodal Machine Learning: A Survey and Taxonomy”. In: *IEEE Transactions on Pattern Analysis and Machine Intelligence* 41.2, pp. 423–443. ISSN: 0162-8828. DOI: [10.1109/TPAMI.2018.2798607](#).
- Caltagirone, Luca et al. (2018). “LIDAR-Camera Fusion for Road Detection Using Fully Convolutional Neural Networks”. In: *arXiv e-prints*, arXiv:1809.07941, arXiv:1809.07941. arXiv: [1809.07941 \[cs.CV\]](#).
- Cayton, Lawrence (2005). “Algorithms for manifold learning”. In: Chauvin, Yves and David E. Rumelhart, eds. (1995). *Backpropagation: Theory, Architectures, and Applications*. Hillsdale, NJ, USA: L. Erlbaum Associates Inc. ISBN: 0-8058-1259-8.
- Cocktail party effect (2010). *Cocktail party effect — Wikipedia, The Free Encyclopedia*. [Online; accessed 29-April-2019]. URL: https://en.wikipedia.org/wiki/Cocktail_party_effect.
- Desimone, Robert and John Duncan (1995). “Neural Mechanisms of Selective Visual Attention”. In: *Annual Review of Neuroscience* 18.1. PMID: 7605061, pp. 193–222. DOI: [10.1146/annurev.ne.18.030195.001205](#). eprint: <https://doi.org/10.1146/annurev.ne.18.030195.001205>. URL: <https://doi.org/10.1146/annurev.ne.18.030195.001205>.
- D.R., Lorimer and Kramer M. (2005). *Handbook of pulsar astronomy*. Cambridge University Press.
- Driver, Jon and Charles Spence (1998). “Crossmodal attention”. In: *Current Opinion in Neurobiology* 8.2, pp. 245 –253. ISSN: 0959-4388. DOI: [https://doi.org/10.1016/S0959-4388\(98\)80147-5](#). URL: <http://www.sciencedirect.com/science/article/pii/S0959438898801475>.
- Ephrat, Ariel et al. (2018). “Looking to Listen at the Cocktail Party: A Speaker-Independent Audio-Visual Model for Speech Separation”. In: *arXiv e-prints*, arXiv:1804.03619, arXiv:1804.03619. arXiv: [1804.03619 \[cs.SD\]](#).
- Fan, Jianqing, Cong Ma, and Yiqiao Zhong (2019). “A Selective Overview of Deep Learning”. In: *arXiv e-prints*, arXiv:1904.05526, arXiv:1904.05526. arXiv: [1904.05526 \[stat.ML\]](#).
- Fazekas, Peter and Bence Nanay (Oct. 2018). “Attention Is Amplification, Not Selection”. In: *The British Journal for the Philosophy of Science*. ISSN: 0007-0882. DOI: [10.1093/bjps/axy065](#). eprint: <http://oup.prod.sis.lan/bjps/advance-10.1093/bjps/axy065>.

- [article-pdf/doi/10.1093/bjps/axy065/25875820/axy065.pdf](https://doi.org/10.1093/bjps/axy065.pdf). URL: <https://doi.org/10.1093/bjps/axy065>.
- Galassi, Andrea, Marco Lippi, and Paolo Torroni (2019). “Attention, please! A Critical Review of Neural Attention Models in Natural Language Processing”. In: *arXiv e-prints*, arXiv:1902.02181, arXiv:1902.02181. arXiv: [1902.02181 \[cs.CL\]](https://arxiv.org/abs/1902.02181).
- Ghahramani, Z. (2004). “Unsupervised Learning”. In: *Springer*.
- Ghosh, Pranab (2007). “Rotation and Accretion Powered Pulsars”. In: *World Scientific Series in Astronomy and Astrophysics*.
- Goodfellow, Ian, Yoshua Bengio, and Aaron Courville (2016). *Deep Learning*. <http://www.deeplearningbook.org>. MIT Press.
- Goodfellow, Ian et al. (2014). “Generative Adversarial Nets”. In: *Advances in Neural Information Processing Systems 27*. Ed. by Z. Ghahramani et al. Curran Associates, Inc., pp. 2672–2680. URL: <http://papers.nips.cc/paper/5423-generative-adversarial-nets.pdf>.
- He, K. et al. (2016). “Deep Residual Learning for Image Recognition”. In: *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 770–778. DOI: [10.1109/CVPR.2016.90](https://doi.org/10.1109/CVPR.2016.90).
- Hoogi, Assaf et al. (2019). “Self-Attention Capsule Networks for Image Classification”. In: *arXiv e-prints*, arXiv:1904.12483, arXiv:1904.12483. arXiv: [1904.12483 \[cs.CV\]](https://arxiv.org/abs/1904.12483).
- Kahneman, Daniel (1975). “Attention and effort”. In:
- Kamyshanska, Hanna and Roland Memisevic (2014). “The Potential Energy of an Autoencoder”. In: *IEEE TRANSACTIONS ON PATTERN ANALYSIS AND MACHINE INTELLIGENCE (TPAMI)*.
- Kim, Taesup and Yoshua Bengio (2016). “Deep Directed Generative Models with Energy-Based Probability Estimation”. In: *arXiv e-prints*, arXiv:1606.03439, arXiv:1606.03439. arXiv: [1606.03439 \[cs.LG\]](https://arxiv.org/abs/1606.03439).
- Kingma, Diederik P. and Jimmy Ba (2014). “Adam: A Method for Stochastic Optimization”. In: *arXiv e-prints*, arXiv:1412.6980, arXiv:1412.6980. arXiv: [1412.6980 \[cs.LG\]](https://arxiv.org/abs/1412.6980).
- Kingma, Diederik P and Max Welling (2013). “Auto-Encoding Variational Bayes”. In: *arXiv e-prints*, arXiv:1312.6114, arXiv:1312.6114. arXiv: [1312.6114 \[stat.ML\]](https://arxiv.org/abs/1312.6114).
- Ladjal, Saïd, Alasdair Newson, and Chi-Hieu Pham (2019). “A PCA-like Autoencoder”. In: *arXiv e-prints*, arXiv:1904.01277, arXiv:1904.01277. arXiv: [1904.01277 \[cs.CV\]](https://arxiv.org/abs/1904.01277).
- LeCun, Yann, Yoshua Bengio, and Geoffrey Hinton (May 2015). “Deep learning”. English (US). In: *Nature* 521.7553, pp. 436–444. ISSN: 0028-0836. DOI: [10.1038/nature14539](https://doi.org/10.1038/nature14539).
- LeCun, Yann et al. (2006). “A tutorial on energy-based learning”. In: *PREDICTING STRUCTURED DATA*. MIT Press.
- Li, Guanbin et al. (2019). “Cross-Modal Attentional Context Learning for RGB-D Object Detection”. In: *IEEE Transactions on Image Processing* 28.4, pp. 1591–1601. DOI: [10.1109/TIP.2018.2878956](https://doi.org/10.1109/TIP.2018.2878956). arXiv: [1810.12829 \[cs.CV\]](https://arxiv.org/abs/1810.12829).
- Li, Yuxi (2017). “Deep Reinforcement Learning: An Overview”. In: *arXiv e-prints*, arXiv:1701.07274, arXiv:1701.07274. arXiv: [1701.07274 \[cs.LG\]](https://arxiv.org/abs/1701.07274).
- Libovický, Jindřich, Jindřich Helcl, and David Mareček (2018). “Input Combination Strategies for Multi-Source Transformer Decoder”. In: *arXiv e-prints*, arXiv:1811.04716, arXiv:1811.04716. arXiv: [1811.04716 \[cs.CL\]](https://arxiv.org/abs/1811.04716).
- Loog, Marco (2017). “Supervised Classification: Quite a Brief Overview”. In: *arXiv e-prints*, arXiv:1710.09230, arXiv:1710.09230. arXiv: [1710.09230 \[cs.LG\]](https://arxiv.org/abs/1710.09230).
- Lorimer, R. Duncan (2008). “Binary and Millisecond Pulsars”. In: *Living Reviews in Relativity* 11.1, p. 8. ISSN: 1433-8351. DOI: [10.12942/lrr-2008-8](https://doi.org/10.12942/lrr-2008-8). URL: <https://doi.org/10.12942/lrr-2008-8>.

- Lyon, R. J. (2016). "Why are pulsars hard to find". PhD thesis. The University of Manchester.
- Narayanan, Hariharan and Sanjoy Mitter (2010). "Sample Complexity of Testing the Manifold Hypothesis". In: *Advances in Neural Information Processing Systems 23*. Ed. by J. D. Lafferty et al. Curran Associates, Inc., pp. 1786–1794. URL: <http://papers.nips.cc/paper/3958-sample-complexity-of-testing-the-manifold-hypothesis.pdf>.
- Ruder, Sebastian (2016). "An overview of gradient descent optimization algorithms". In: *arXiv e-prints*, arXiv:1609.04747, arXiv:1609.04747. arXiv: [1609.04747 \[cs.LG\]](#).
- Santilli, RM (1982). "Birkhoffian generalization of Hamiltonian Mechanics". In: *Foundations of theoretical mechanics II*.
- Scholz, Matthias, Martin Fraunholz, and Joachim Selbig (2008). "Nonlinear Principal Component Analysis: Neural Network Models and Applications". In: *Principal Manifolds for Data Visualization and Dimension Reduction*. Ed. by Alexander N. Gorban et al. Berlin, Heidelberg: Springer Berlin Heidelberg, pp. 44–67. ISBN: 978-3-540-73750-6.
- Shon, Suwon, Tae-Hyun Oh, and James Glass (2018). "Noise-tolerant Audio-visual Online Person Verification using an Attention-based Neural Network Fusion". In: *arXiv e-prints*, arXiv:1811.10813, arXiv:1811.10813. arXiv: [1811.10813 \[cs.CV\]](#).
- Treat, J. and Stegmaier (2014). "Black Holes: Star Eater." In: *National Geographic*.
- Vaswani, Ashish et al. (2017). "Attention Is All You Need". In: *arXiv e-prints*, arXiv:1706.03762, arXiv:1706.03762. arXiv: [1706.03762 \[cs.CL\]](#).
- Vincent, Pascal et al. (2008). "Extracting and Composing Robust Features with Denoising Autoencoders". In: *ICML 2008*.
- Wang, Xin, Yuan-Fang Wang, and William Yang Wang (2018). "Watch, Listen, and Describe: Globally and Locally Aligned Cross-Modal Attentions for Video Captioning". In: *arXiv e-prints*, arXiv:1804.05448, arXiv:1804.05448. arXiv: [1804.05448 \[cs.CL\]](#).
- Watzl, Sebastian (2017). *Structuring Mind. The Nature of Attention and How It Shapes Consciousness*. Oxford, UK: Oxford University Press.
- Weinberger, K.Q. and L.K. Saul (2006). "Unsupervised Learning of Image Manifolds by Semidefinite Programming". In: *Int J Comput Vision*.
- Wu, Yonghui et al. (2016). "Google's Neural Machine Translation System: Bridging the Gap between Human and Machine Translation". In: *ArXiv* abs/1609.08144.
- Zhai, Shuangfei et al. (2016). "Deep Structured Energy Based Models for Anomaly Detection". In: *arXiv e-prints*, arXiv:1605.07717, arXiv:1605.07717. arXiv: [1605.07717 \[cs.LG\]](#).